

---

# Damped Anderson Mixing for Deep Reinforcement Learning: Acceleration, Convergence, and Stabilization

---

Ke Sun<sup>\*1</sup>, Yafei Wang<sup>\*1</sup>, Yi Liu<sup>1</sup>, Yingnan Zhao<sup>1,2</sup>, Bo Pan<sup>1</sup>, Shangling Jui<sup>3</sup>,  
Bei Jiang<sup>1</sup>, Linglong Kong<sup>1†</sup>

<sup>1</sup>University of Alberta, Edmonton, Canada

<sup>2</sup>Harbin Institute of Technology, Harbin, China

<sup>3</sup>Huawei Technologies Ltd.

{ksun6,yafei2,yliu16,yingnan6,pan1,bei1,lkong}@ualberta.ca  
jui.shangling@huawei.com

## Abstract

Anderson mixing has been heuristically applied to reinforcement learning (RL) algorithms for accelerating convergence and improving the sampling efficiency of deep RL. Despite its heuristic improvement of convergence, a rigorous mathematical justification for the benefits of Anderson mixing in RL has not yet been put forward. In this paper, we provide deeper insights into a class of acceleration schemes built on Anderson mixing that improve the convergence of deep RL algorithms. Our main results establish a connection between Anderson mixing and quasi-Newton methods and prove that Anderson mixing increases the convergence radius of policy iteration schemes by an extra contraction factor. The key focus of the analysis roots in the fixed-point iteration nature of RL. We further propose a stabilization strategy by introducing a stable regularization term in Anderson mixing and a differentiable, non-expansive MellowMax operator that can allow both faster convergence and more stable behavior. Extensive experiments demonstrate that our proposed method enhances the convergence, stability, and performance of RL algorithms.

## 1 Introduction

In reinforcement learning (RL) [1], an agent seeks an optimal policy in a sequential decision-making process. Deep RL has recently achieved significant improvements in a variety of challenging tasks, including game playing [2, 3, 4] and robust navigation [5]. A flurry of state-of-the-art algorithms have been proposed, including Deep Q-Learning (DQN) [2] and variants such as Double-DQN [6], Dueling-DQN [7], Deep Deterministic Policy Gradient (DDPG) [8], Soft Actor-Critic [9] and distributional RL algorithms [10, 11, 12], all of which have successfully solved end-to-end decision-making problems such as playing Atari games. However, the slow convergence and sample inefficiency of RL algorithms still hinders the progress of RL research, particularly in high-dimensional state spaces where deep neural network are used as function approximators, making learning in real physical worlds impractical.

To address these issues, various acceleration strategies have been proposed, including the classical Gauss-Seidel Value Iteration [13] and Jacobi Value Iteration [14]. Another popular branch of

---

\*Equal contributions in alphabetical order

†Corresponding author

techniques accelerates RL by leveraging historical data. Interpolation methods such as Average-DQN [15] have been widely used in first-order optimization problems [16] and have been proven to converge faster than vanilla gradient methods. As an effective multi-step interpolation method, Anderson mixing [17, 18], also known as Anderson acceleration, has attracted great attention from RL researchers. The insight underpinning of Anderson acceleration is that RL [1] is intrinsically linked to fixed-point iterations: the optimal value function is the fixed point of the Bellman optimality operator. These fixed-points are computed recursively by repeatedly applying an operator of interest [19]. Anderson mixing is a general method to accelerate fixed-point iterations [17] and has been successfully applied to fields, such as the computational chemistry [20] or electronic structure computation [21]. In particular, Anderson acceleration leverages the  $m$  previous estimates in order to find a better estimate in a fixed-point iteration. To compute the mixing coefficients in Anderson iteration, it searches for a point with a minimal residual within the subspace spanned by these estimates. It is thus natural to explore the efficacy of Anderson acceleration in RL settings.

Several works [22, 23] have attempted to apply Anderson acceleration to reinforcement learning. Anderson mixing was first applied to value iteration in [19, 22] and resulted in significant convergence improvements. Regularized Anderson acceleration [23] was recently proposed to further accelerate convergence and enhance the final performance of state-of-the-art RL algorithms in various experiments. However, previous applications of Anderson acceleration were typically heuristic: consequently, these empirical improvements in convergence have so far lacked a rigorous mathematical justification.

In this paper, we provide deeper insights into Anderson acceleration in reinforcement learning by establishing its connection with quasi-Newton methods for policy iteration and improved convergence guarantees under the assumptions that the Bellman operator is differential and non-expansive. MellowMax operator is adopted to replace the max operator in policy iteration to simultaneously guarantee faster convergence of value function and reduce the estimated gradient variance to yield stabilization. In addition, we analyze the stability properties of Anderson acceleration in policy iteration and propose a stable regularization to further enhance the stability. These key two factors, i.e., the stable regularization and the theoretically-inspired MellowMax operator, are the basis for our *Stable Anderson Acceleration (Stable AA)* method. Finally, our experimental results on various Atari games demonstrate that our Stable AA method enjoys faster convergence and achieves better performance relative to existing Anderson acceleration baselines. Our work provides a unified analytic framework that illuminates Anderson acceleration for reinforcement learning algorithms from the perspectives of acceleration, convergence, and stabilization.

## 2 Acceleration and Convergence Analysis of Anderson Acceleration on RL

We first present the notion of Anderson acceleration in the reinforcement learning and then provide deeper insights into the acceleration it affords by establishing a connection with quasi-Newton methods. Finally, a theoretical convergence analysis is provided to demonstrate that Anderson acceleration can increase the convergence radius of policy iteration by an extra contraction factor.

**Background** Consider a Markov decision process (MDP) specified by the tuple  $(\mathcal{S}, \mathcal{A}, R, p, \gamma)$ , where  $\mathcal{S}$  is a set of the states and  $\mathcal{A}$  is a set of actions. The functions  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  and  $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  are the reward function, with  $R_t = R(s, a)$ , and transition dynamics function, respectively for the MDP. The discount rate is denoted by  $\gamma \in [0, 1)$  and determines the relative importance of immediate rewards relative to rewards received in the future. The  $Q$ -value function evaluates the expected return starting from a given state-action pair  $(s, a)$ , that is,  $Q^\pi(s, a) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid s_0 = s, a_0 = a]$ . A policy  $\pi(a|s)$  is a distribution mapping the state space  $\mathcal{S}$  to the action space  $\mathcal{A}$ .

### 2.1 Anderson Acceleration in Policy Iteration

We focus on the tabular case to enable the theoretical analysis of Anderson acceleration in value (policy) iteration, which can be naturally applied to function approximation. Both the value iteration ( $V$ -notation) and the policy iteration ( $Q$ -notation) can have Anderson acceleration applied to them to improve convergence. However, theoretical analysis has shown that value iteration enjoys a  $\gamma$ -linear convergence rate, i.e.,  $\|V^{(t)} - V^*\|_\infty \leq \gamma \|V^{(t-1)} - V^*\|_\infty$ , where  $V^{(t)}$  is the value

function in the iteration step  $t$  and  $V^*$  is the optimal value function, while policy iteration converges faster. This is due to the fact that policy iteration more fully evaluates the current policy than does value iteration. Additionally, policy iteration is more fundamental and scales more readily to deep reinforcement learning. For this reason, we analyze our method in the policy iteration setting under  $Q$ -notation as value iteration is a special case of policy iteration. Thus, our analysis also applies under value iteration.

We first focus on the control setting where the optimal value of state-action pair  $Q^*(s, a)$  is defined recursively as a function of the optimal value of the other state-action pair:

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) \cdot \max_{a'} Q^*(s', a'). \quad (1)$$

Combining policy evaluation and policy improvement, the resulting policy iteration algorithm is equivalent to solving for the fixed point of the Bellman optimality operator  $\mathcal{T} : \mathbb{R}^{|S \times A|} \rightarrow \mathbb{R}^{|S \times A|}$  with  $\mathcal{T}Q(s, a) = R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) \cdot \max_{a'} Q(s', a')$ .

As a general technique to speed up fixed-point iteration [17], Anderson acceleration has been successfully yet heuristically applied to reinforcement learning algorithms [23, 19, 22]. Specifically, Anderson acceleration linearly combines the  $m$  previous estimates to yield a better iteration target in the fixed point iteration. Geometrically, Anderson acceleration applies the operator to a point that has a minimal residual within the subspace spanned by these estimates. In policy iteration, Anderson acceleration maintains a memory of the previous  $Q$  function values and updates the iterate as a linear combination of these values with dynamic weights  $\alpha^{(k)}$  in the  $k$ th iteration step. Specifically,

$$Q^{(k+1)}(s, a) = (1 - \beta_k) \sum_{i=0}^m \alpha_i^{(k)} Q^{(k-m+i)}(s, a) + \beta_k \sum_{i=0}^m \alpha_i^{(k)} \mathcal{T}Q^{(k-m+i)}(s, a), \quad (2)$$

where  $0 \leq \beta_k \leq 1$  is the damping parameter. All of the coefficients  $\alpha_i^{(k)}$  in the coefficient vector  $\alpha^{(k)}$  are computed following

$$\alpha^{(k)} = \underset{\alpha \in \mathbb{R}^{m+1}}{\operatorname{argmin}} \left\| \sum_{i=0}^m \alpha_i \left( \mathcal{T}Q^{(k-m+i)} - Q^{(k-m+i)} \right) \right\|_2 = \underset{\alpha \in \mathbb{R}^{m+1}}{\operatorname{argmin}} \left\| \Delta_k'^T \cdot \alpha \right\|_2, \text{ s.t. } \sum_{i=0}^m \alpha_i = 1, \quad (3)$$

where  $\Delta_k'^T = [e_{k-m}, \dots, e_k] \in \mathbb{R}^{|S \times A| \times (m+1)}$  and  $e_k = \mathcal{T}Q^{(k)} - Q^{(k)} \in \mathbb{R}^{|S \times A|}$  is the Bellman residuals matrix. By the Karush-Kuhn-Tucker conditions, the analytic solution of optimal coefficient vector  $\alpha^k$  is

$$\alpha^{(k)} = \frac{\left( \Delta_k' \Delta_k'^T \right)^{-1} \mathbf{1}}{\mathbf{1}^T \left( \Delta_k' \Delta_k'^T \right)^{-1} \mathbf{1}}, \quad (4)$$

where  $\mathbf{1}$  denotes the vector with all components equal to one.

## 2.2 Connection Between Damped Anderson Acceleration and Quasi-Newton Methods

We know that the optimization problem is closely linked with solving a fixed-point iteration problem by directly solving its first-order condition. Inspired by [24], we show that Anderson acceleration in policy iteration attempts to perform a special form of quasi-Newton iteration from its optimization problem behind.

To illuminate this connection, we firstly show that the original constrained optimization to obtain the optimal  $\alpha^{(k)}$  in Eq.(3) can be equivalent to the unconstrained one

$$\tau^{(k)} = \underset{\tau \in \mathbb{R}^m}{\operatorname{argmin}} \|e_k - H_k \tau\|^2, \quad (5)$$

where we let  $n = |S| \times |A|$ , and then  $H_k = [e_k - e_{k-1}, \dots, e_{k-m+1} - e_{k-m}] \in \mathbb{R}^{n \times m}$ .  $\tau^{(k)} = [\tau_0^{(k)}, \tau_1^{(k)}, \dots, \tau_{m-1}^{(k)}]^T \in \mathbb{R}^m$  with  $\tau_i^{(k)} = \sum_{j=0}^{m-i-1} \alpha_j^{(k)}$ . Let  $\delta_k = Q^{(k)} - Q^{(k-1)}$ ,  $\Delta_k = [\delta_k, \delta_{k-1}, \dots, \delta_{k-m+1}] \in \mathbb{R}^{n \times m}$ . We show that the updating rule of  $Q(s, a)$  in Anderson acceleration can be expressed as a quasi-Newton form in Proposition 1.

**Proposition 1.** *By conducting the damped Anderson acceleration (Eq.(2) and (3)) on the policy iteration, the updating of  $Q^{(k+1)}$  can be reformulated as*

$$Q^{(k+1)} := Q^{(k)} - G_k e_k \quad (6)$$

where  $G_k = (\Delta_k + \beta_k H_k) (H_k^T H_k)^{-1} H_k^T - \beta_k I$  can serve as an approximation of inverse Jacobian matrix of  $e_k = TQ^{(k)} - Q^{(k)}$ , and  $I$  is an identity matrix.

Proposition 1 points out that Anderson acceleration on policy iteration additionally leverages more information about the fixed-point residual  $e_k = TQ^{(k)} - Q^{(k)}$  to update the  $Q$  function. Particularly, the first part  $(\Delta_k + \beta_k H_k)(H_k^T H_k)^{-1} H_k^T$  in  $G_k$  contains partial structure matrix information about the real inverse of Jacobian matrix, which has the huge potential to speed up the convergence of the fixed-point iteration. More importantly, the results established in [17] and [18] can be seen as special cases of Proposition 1 with  $\beta_k = 1$ . If we directly get rid of the first part in  $G_k$  and set  $\beta_k = 1$ , the updating rule exactly degenerates to the Q-value function iteration without Anderson acceleration.

### 2.3 Convergence Rate Analysis of Anderson Acceleration on RL

The success of the Anderson acceleration to reduce the residual is coupled in the algorithm iteration at each stage. Let  $e_k^\alpha = \sum_{j=0}^m \alpha_j^{(k)} (TQ^{(j)} - Q^{(j)})$ . The stage- $k$  gain  $\theta_k$  can be defined by  $\|e_k^\alpha\|_\infty = \theta_k \|e_k\|_\infty$ . As  $\alpha_k^{(k)} = 1, \alpha_j^{(k)} = 0, j \neq k$ , i.e.,  $m = 0$ , is an admissible solution to the optimization problem in Eq. (3), it immediately follows that  $0 \leq \theta_k \leq 1$ . The key to rigorously show that Anderson acceleration can speed up the convergence of policy iteration by taking a linear combination of history steps is connecting the gain  $\theta_k$  to the differences of consecutive iterates  $Q^{(k)}$  and residual terms  $e_k$ . As discussed in the following part, the improvement of the convergence rate of the policy iteration by using the acceleration technique is characterized by  $\theta_k$ . We first consider the following assumption about the operator  $\mathcal{T}$  used to guarantee the first and second order derivatives of  $\mathcal{T}$  are bounded, as in [25].

**Assumption 1.** *Assume the Bellman operator  $\mathcal{T}$  acting on state-action value function  $Q$  has a fixed point  $Q^*$ , and there are positive constants  $c_1$  and  $c_2$  such that*

1.  $\mathcal{T} \in C^2(\mathbb{R}^{|\mathcal{S} \times \mathcal{A}|})$ .
2. *The first derivative of  $\mathcal{T}$  is bounded by  $c_1$ .*
3. *The second derivative of  $\mathcal{T}$  is bounded by  $c_2$ .*

**Theorem 1.** *Under Assumption 1, if the coefficients  $\alpha_i^{(k)}$  remain bounded and away from zero, the following bound holds for the fixed point residual  $e_k$  from Eq. (2) and (3) with depth  $m$*

$$\begin{aligned} \|e_k\|_\infty &\leq \theta_k \left\{ (1 - \beta_{k-1}) + c_1 \beta_{k-1} \right\} \|e_{k-1}\|_\infty + c_2 \cdot (\|\delta_k\|_\infty + \|\delta_{k-1}\|_\infty) |\tau_1| \|\delta_{k-1}\|_\infty \\ &+ c_2 \cdot \sum_{i=2}^m \left( \|\delta_k\|_\infty + \|\delta_{k-i}\|_\infty + 2 \sum_{l=1}^{i-1} \|\delta_{k-l}\|_\infty \right) |\tau_i| \|\delta_{k-i}\|_\infty. \end{aligned} \quad (7)$$

Note that the first term of RHS in Eq. (7) characterizes an increased convergence radius by an extra contraction factor  $\theta_k$ . Therefore, if the error terms  $\|\delta_{k-i}\|_\infty$  are small enough, and the operator  $\mathcal{T}$  is differentiable and has bounded first and second derivatives, the faster convergence result characterized by  $\theta_k$  can be derived for the Q-value function iteration with Anderson acceleration.

Unfortunately, the commonly used max operator in  $\mathcal{T}$  does not satisfy Assumption 1 as it is not a differentiable operator. Moreover, the ‘‘hard’’ max operator in  $\mathcal{T}$  always commits to the maximum action-value function according to current estimation for updating the value estimator, lacking the ability to consider other potential action-values. A natural alternative is the Boltzmann Softmax operator, but this operator is prone to misbehave [26] as it is not a non-expansive operator. MellowMax operator [26], which can help strike a balance between exploration and exploitation, is considered in this paper. More importantly, the more meaningful convergence result of Anderson acceleration in policy iteration under Assumption 2 can be established due to the contraction properties of the Bellman operator under MellowMax operator. The results are given in Theorem 2.

**Assumption 2.** Assume the Bellman operator  $\mathcal{T}$  acting on state-action value function  $Q$  is a  $\gamma$ -contraction operator, i.e.,  $\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma\|Q - Q'\|_\infty$  for each state-action function pair  $Q$  and  $Q'$ .

**Theorem 2.** If both Assumption 1 and 2 hold, the coefficients  $\alpha_i^{(k)}$  remain bounded and away from zero. The following bound holds for the residual  $e_k$  with depth  $m$

$$\|e_k\|_\infty \leq \theta_k[(1 - \beta_{k-1}) + c_1\beta_{k-1}]\|e_{k-1}\|_\infty + O\left(\sum_{j=1}^m \|e_{k-j}\|_\infty^2\right). \quad (8)$$

From Theorem 2, we find that there is a theoretical advantage to consider Anderson acceleration for policy iteration with depth  $m$  due to the gain  $\theta_k$  even with the higher-order accumulating terms. Fortunately, the Bellman operator with MellowMax operator is a  $\gamma$ -contraction operator, satisfying Assumption 2. Specifically, the resulting Bellman operator  $\mathcal{T}_{mm}$  under the MellowMax operator is defined as

$$\mathcal{T}_{mm}Q(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) mm_\omega(Q(s', \cdot)), \quad (9)$$

where  $mm_\omega$  is the MellowMax operator and  $mm_\omega Q(s', \cdot) = \log(\frac{1}{|\mathcal{A}|} \sum_{a'} \exp[\omega Q(s', a')])/\omega$ . Rigorously, we show that MellowMax can simultaneously satisfy Assumption 1 and 2 in Appendix B. As such, the faster convergence result presented in Theorem 2 can be derived. In other words, the faster convergence of policy iteration under MellowMax operator can be established by applying Anderson acceleration. Based on this theoretical principle, we apply MelloxMax operator in the Bellman operator to design our method in Section 3, where a detailed discussion is also provided.

Finally, the  $Q$ -value estimate  $Q^{(k)}$  can be obtained by iteratively applying the MellowMax operator by starting from some initial value  $Q^{(0)}$ :

$$Q^{(k+1)} \leftarrow (1 - \beta_k) \sum_{i=0}^m \alpha_i^{(k)} Q^{(k-m+i)} + \beta_k \sum_{i=0}^m \alpha_i^{(k)} \mathcal{T}_{mm}Q^{(k-m+i)}, \quad \forall (s, a) \in (\mathcal{S}, \mathcal{A}). \quad (10)$$

### 3 Stabilization Analysis and Our Method

In this section, a stable regularization is firstly introduced and its stability analysis is provided as well. We then briefly analyze the role that MellowMax operator plays when conducting Anderson acceleration on deep reinforcement learning. These two factors eventually inspire our algorithm, which we call *Stable Anderson Acceleration (Stable AA)*.

#### 3.1 Stable Regularization

Inspired by recent stable results of Anderson acceleration [18], we introduce the stable regularization term on the aforementioned unconstrained optimization problem Eq. (5) to obtain mixing coefficients  $\tau^k$ . Particularly, we add  $\ell_2$  regularization of  $\tau^k$  scaled by the Frobenius norm of  $\Delta_k$  and  $H_k$  to improve the stability. This yields the new optimization problem

$$\tau^k = \underset{\tau \in \mathbb{R}^m}{\operatorname{argmin}} \|e_k - H_k \tau\|^2 + \eta \left( \|\Delta_k\|_F^2 + \|H_k\|_F^2 \right) \|\tau\|^2, \quad (11)$$

where  $\eta$  is a positive tuning parameter representing the scale of regularization. The solution is  $\tau^k = (H_k^T H_k + \eta(\|\Delta_k\|_F^2 + \|H_k\|_F^2)\mathbf{I})^{-1} H_k^T e_k$ . We introduce this stable regularization under the unconstrained variables  $\tau^k$ , which facilitates the optimization. Intuitively, if the algorithm converges, we have  $\lim_{k \rightarrow \infty} \|\Delta_k\|_F = \lim_{k \rightarrow \infty} \|H_k\|_F = 0$ . Therefore, the coefficient on the regularization term vanishes as the algorithm converges, degenerating to Anderson acceleration method without the stable regularization. In this sense, the stability owing to our stable regularization plays a more important role in the early phase of training, which we demonstrate in Section 4. Based on the solved stable regularization  $\tau^k$  and the relationship between  $\tau^k$  and  $\alpha^{(k)}$ , we derive the updating of  $Q^{(k+1)}$  in the policy iteration as follows

$$Q^{(k+1)} = Q^{(k)} - \tilde{G}_k e_k, \quad (12)$$

where  $\tilde{G}_k = -\beta_k \mathbf{I} + (\Delta_k + \beta_k H_k) (H_k^T H_k + \eta(\|\Delta_k\|_F^2 + \|H_k\|_F^2)\mathbf{I})^{-1} H_k^T$ .

Moreover, the following Theorem characterizes the stability ensured by regularization in Eq (11). Please refer to the proof in Appendix C.

**Theorem 3.** *The matrix  $\tilde{G}_k$  satisfy  $\|\tilde{G}_k\|_2 \leq |2/\eta - \beta_k|$ ,  $\|\tilde{G}_k^{-1}G_k\|_2 < 1$ .*

Theorem 3 derives the upper bound of  $\|\tilde{G}_k\|_2$ , which is determined by  $\eta$  and  $\beta_k$ . Intuitively, a larger strength of regularization  $\eta$  and a proper magnitude of  $\beta_k$  can yield more stability. In addition,  $\|\tilde{G}_k^{-1}G_k\|_2$  is strictly less than 1, revealing a smaller violation in  $Q^{(k)}$  iteration compared with non-regularized form.

To quantify the effect of regularization on the coefficient  $\alpha^{(k)}$ , we provide some analytical results regarding the obtained mixing coefficients  $\alpha^{(k)}$  in Proposition 2. The proof is provided in Appendix C.

**Proposition 2.** *Let  $\alpha_{non}^{(k)}$  and  $\alpha_{reg}^{(k)}$  be the mixing coefficient vectors obtained by vanilla unconstrained and our stable regularized Anderson acceleration, respectively. Define the transformation matrix as  $A$ , satisfying  $\alpha^{(k)} = A \cdot \tilde{\tau}^k$  with  $\tilde{\tau}^{(k)} = (1, \tau^k)^T$  (detailed structure of  $A$  is in the Appendix C). Let  $cond_2(A)$  be the conditional number of  $A$ , i.e.,  $cond_2(A) = \|A\|_2 \|A^{-1}\|_2$ . Then we have the following inequalities*

$$\|\alpha_{reg}^{(k)}\|_2^2 \leq 4\left(1 + \frac{\|e_k\|^2}{\eta^2}\right), \quad \|\alpha_{reg}^{(k)} - \alpha_{non}^{(k)}\|_2^2 \leq (cond_2(A))^2 \cdot \left\|\alpha_{non}^{(k)}\right\|_2^2 - \frac{2m+1}{m+1}. \quad (13)$$

From the first inequality, we observe that the  $\ell_2$ -norm of the derived coefficients  $\alpha_{reg}^{(k)}$  is controlled by the regularization parameter  $\eta$ . An overly large  $\eta$  tends to reduce the bound for the norm of  $\alpha_{reg}^{(k)}$ , implying a stable solution of the mixing coefficients  $\alpha_{reg}^{(k)}$ . Besides, we can conclude from the second inequality that there is an inevitable gap between  $\alpha_{non}^{(k)}$  and  $\alpha_{reg}^{(k)}$ .

### 3.2 Stability Effects of MellowMax

The adopted MellowMax operator bridges the Anderson acceleration and reinforcement learning algorithms and it has two-sided stability effects. Firstly, based on the convergence analysis in Section 2, MellowMax operator satisfies the *differential and non-expansive* properties, which allows the faster convergence of Anderson acceleration in policy iteration. In contrast, the commonly used max and Boltzmann Softmax operator [27, 1] violate one of the theoretical assumptions respectively, and thus the (faster) convergence of Anderson acceleration under them may not be guaranteed. This is likely to yield instability while the training of algorithms.

Secondly, it is well-known that the ‘‘hard’’ max updating scheme in the popular off-policy methods, such as Q-learning [28], may lead to misbehavior due to the overestimation issue in the noisy environment [29, 6, 30]. By contrast, it has been demonstrated that MellowMax and Softmax operators are capable of reducing the overestimation biases, therefore reducing the gradient noises to stabilize the optimization of neural networks [26, 31]. The stable gradient estimation based on MellowMax operator leads to the enhancement of final performance for algorithms.

### 3.3 Algorithm: Stable AA

The introduced stable regularization approach combined with the MellowMax operator finally form our Stable AA method. In our algorithm, we focus on exploring the impact of Stable AA on Dueling-DQN [7]. In particular, under the procedure of off-policy learning in DQN, we firstly formulate the general damped Anderson acceleration form with the function approximator  $Q_\theta$  as follows

$$Q_\theta(s_t, a_t) = \beta_t \sum_{i=1}^m \hat{\alpha}_i Q_{\theta^i}(s_t, a_t) + (1 - \beta_t) \mathbb{E}_{s_{t+1}, r_t} \left[ r_t + \gamma \sum_{i=1}^m \hat{\alpha}_i \max_{a_{t+1}} Q_{\theta^i}(s_{t+1}, a_{t+1}) \right], \quad (14)$$

where  $\theta_i$  are parameters of target network before the  $i$ -th update step.  $\hat{\alpha}_i$  can be computed either by vanilla Anderson acceleration [19], or Regularized Anderson acceleration [23]. Then the obtained  $Q_\theta(s_t, a_t)$  serves as the target in the updating of Q-networks. In our Stable AA method, we firstly solve the optimization problem in Eq. (11) to compute  $\tau^k$ . Next we obtain  $\tilde{\alpha}^{(k)}$  by making use of the quantitative relationship between  $\tau^k$  and  $\alpha^{(k)}$ . More importantly, we substitute max with

---

**Algorithm 1** Stable AA Dueling-DQN Algorithm

---

- 1: Initialize a Q value network  $Q_\theta$  and  $m$  target networks with parameters  $\theta_i$  ( $i = 1, \dots, m$ ). Set the total training steps  $K$  and updating step  $M$ .
  - 2: **while**  $k \leq K$  **do**
  - 3:   Observe the initial state  $s_0$ ;
  - 4:   **for**  $t = 1$  to  $T$  **do**
  - 5:     Select  $a_t = \arg \max_a Q_\theta(s_t, a)$  with probability  $1 - \epsilon$  and a random action with probability  $\epsilon$ .
  - 6:     Perform the action  $a_t$ , obtain  $r_t$  and  $s_{t+1}$ . Store the transition  $(s_t, a_t, r_t, s_{t+1})$  in the replay buffer.
  - 7:     Sample the batch of transitions  $(s, a, r, s')$  from the replay buffer.
  - 8:     */\* Step 1: compute  $\tilde{\alpha}^{(k)}$  \*/*
  - 9:     Compute  $\Delta_k$  and  $H_k$ , and then solve the optimization problem with stable regularization in Eq. (11) to obtain  $\tau^k$ .
  - 10:    Obtain the optimal coefficient vectors  $\tilde{\alpha}^{(k)}$  via  $\tilde{\alpha}^{(k)} = A \cdot \tilde{\tau}^k$ , where the transformation matrix  $A$  is defined in Proposition 2.
  - 11:    */\* Step 2: compute the target  $\tilde{Q}_\theta$  by Anderson Mixing \*/*
  - 12:    Compute the value after the MellowMax operator for each target network  $Q_{\theta_i}$ , i.e.,  $mm_\omega(Q_{\theta_i}(s_{t+1}, \cdot))$
  - 13:    Evaluate the target value function  $\tilde{Q}_\theta(s_t, a_t)$  via Eq. (15) under the Melloxmax operator.
  - 14:    */\* Step 3: update the Q value networks \*/*
  - 15:    Update the Q value network  $\theta$  by minimizing the loss in Eq.(16) with the target  $y_t$  from Step 2.
  - 16:    Update  $m$  target networks every  $M$  steps, i.e.,  $\theta_i \leftarrow \theta_{i+1}$  ( $i = 1, \dots, m$ ) and  $\theta_m \leftarrow \theta$ .
  - 17:    Set  $k = k + 1$ .
  - 18:    **end for**
  - 19: **end while**
- 

MellowMax operator  $mm_\omega$ . The resulting target value function  $\tilde{Q}_\theta$  in our Stable AA algorithm is reformulated as

$$\tilde{Q}_\theta(s_t, a_t) = \beta_t \sum_{i=1}^m \tilde{\alpha}_i Q_{\theta_i}(s_t, a_t) + (1 - \beta_t) \mathbb{E}_{s_{t+1}, r_t} \left[ r_t + \gamma \sum_{i=1}^m \tilde{\alpha}_i \cdot mm_\omega(Q_{\theta_i}(s_{t+1}, \cdot)) \right], \quad (15)$$

where  $mm_\omega(Q(s, \cdot))$  is the MellowMax operator. Finally, the updating is performed by minimizing the least squared errors of Bellman equation between the current Q value estimate  $Q_\theta(s_t, a_t)$  and the target value function  $y_t$  obtained from Eq. (15),

$$L(\theta) = \mathbb{E}_{(s_t, a_t) \in \mathcal{D}} \left[ (y_t - Q_\theta(s_t, a_t))^2 \right], \quad (16)$$

where  $\mathcal{D}$  is the distribution of previously sampled transitions.

In summary, the key of Stable AA method in policy iteration lies in two factors: the stable regularization in Eq. (11) in computing coefficient  $\alpha^{(k)}$ , and the MellowMax operator enabling the faster convergence in updating  $Q^{(k)}$ , both of which improve the convergence and sample efficiency. Moreover, we provide a detailed description of Stable AA on Dueling-DQN algorithm in Algorithm 1. Similar to the strategy in [23], the incorporation of Stable AA into policy gradient based algorithms, including actor critic [1] and twin delayed DDPG (TD3) [32] can be directly implemented in their critics part. It can be viewed as a straightforward extension, and we leave this exploration as future works.

## 4 Experiment

Our theoretical results about Anderson acceleration mainly apply to the case of a tabular value function representation, but our derived Stable AA algorithms can be naturally applied into the function approximation setting. The goal of our experiments is to show that our Stable AA method can still be useful in practice by improving the performance of Dueling-DQN algorithms. Our experimental results demonstrate that such an improvement is attributed to the joint benefits of the proposed stable regularization and the MelloMax operator.

**Experimental Settings** We perform our Stable AA Dueling-DQN algorithm on a variety of Atari games, and mainly focus on reporting four representative games, i.e., SpaceInvaders, Enduro, Breakout, and CrazyClimber. Results of other games are similar, which we provide in Appendix D. We compare our approach with Dueling-DQN [7] and Regularized Anderson Acceleration (RAA)[23]. In addition, we also provide an ablative analysis about our Stable AA algorithm to illuminate the separate and joint impacts of stable regularization and MellowMax operator. In the following experiments, we report results statistics by running three independent random seeds. We set  $\beta_t$  in Eq. (15) as 0.1 for convenience.

**Implementation of MellowMax Operator** The MellowMax operator  $mm_\omega$  satisfies the desirable differential and non-expansive properties, enabling the faster convergence in policy iteration with Anderson acceleration. Nevertheless, we need to perform an additional root-finding algorithm [26] to compute the optimal  $\omega$  in each state in order to maintain these properties and help the MellowMax operator to identify a probability distribution. Unfortunately, this root-finding algorithm is computationally expensive to be applied. Following the strategy in [31], we tune the inverse parameter  $\omega$  from  $\{1, 5, 10\}$  and then report the best score.

#### 4.1 Performance of Our Stable AA

We select DuelingDQN and DuelingDQN-RAA as baselines for the evaluation on the four Atari games. These two algorithms and our approach DuelingDQN-Stable AA are trained under the same random seeds and evaluated every 10,000 environment steps, where each evaluation reports the average returns with standard deviations. Our implementation is adapted from RAA [23]. After the grid search, we set  $\omega$  in MelloxMax of Stable AA as 5.0 and  $\eta$  in stable regularization as 0.1 across 4 Atari games.

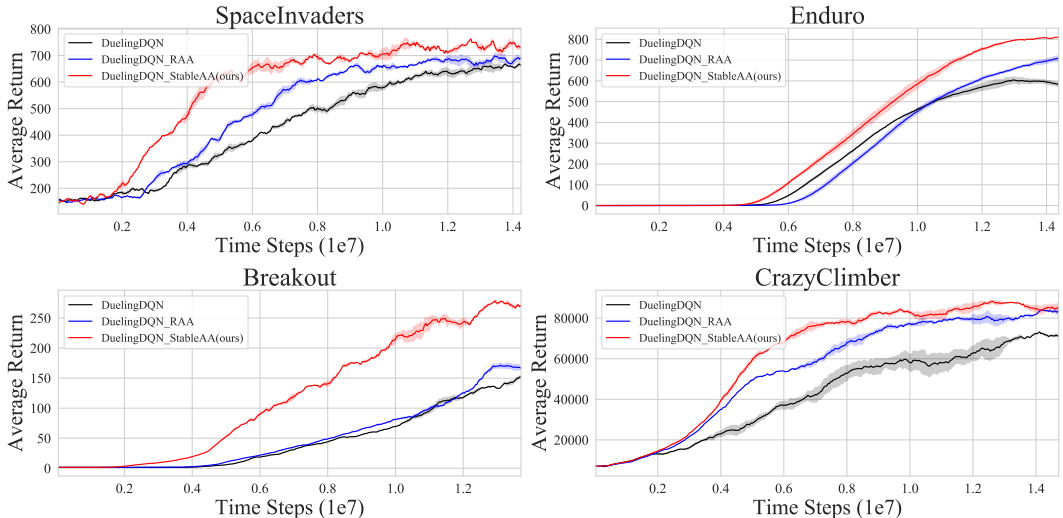


Figure 1: Learning curves of DuelingDQN, DuelingDQN-RAA and our approach DuelingDQN-Stable AA on SpaceInvaders, Enduro, Breakout, and CrazyClimber games over 3 seeds. Shaded region corresponds to the standard deviation.

From Figure 1, we note that our DuelingDQN-StableAA (red line) significantly outperforms Regularized AA (blue line) and baseline (black line) across all four games. Overall, Dueling-RAA enables to accelerate DuelingDQN to improve the sample efficiency and enhance the final performance, but our approach can lead to further benefits. Remarkably, our DuelingDQN-StableAA (red line) is superior to RAA to a large margin, especially on Breakout where our approach achieves around 250 average return while RAA only achieves 150 return. In summary, we conclude that the joint impact of both stable regularization and theoretically-principled MellowMax further accelerate the convergence and improve the sample efficiency of the popular off-policy DuelingDQN algorithm.



## 4.2 Ablation Analysis

We further examine the separate impact of the proposed stable regularization (shown in Eq. (11)) and the theoretically-principled MelloxMax operator via the rigorous ablation study. Starting from DuelingDQN, we firstly add stable regularization with different scales  $\eta$  while comparing with our resulting Stable AA method. Meanwhile, we separately replace Max operator in DuelingDQN with MelloxMax operators with various inverse parameters  $\omega$  to explore their impacts.

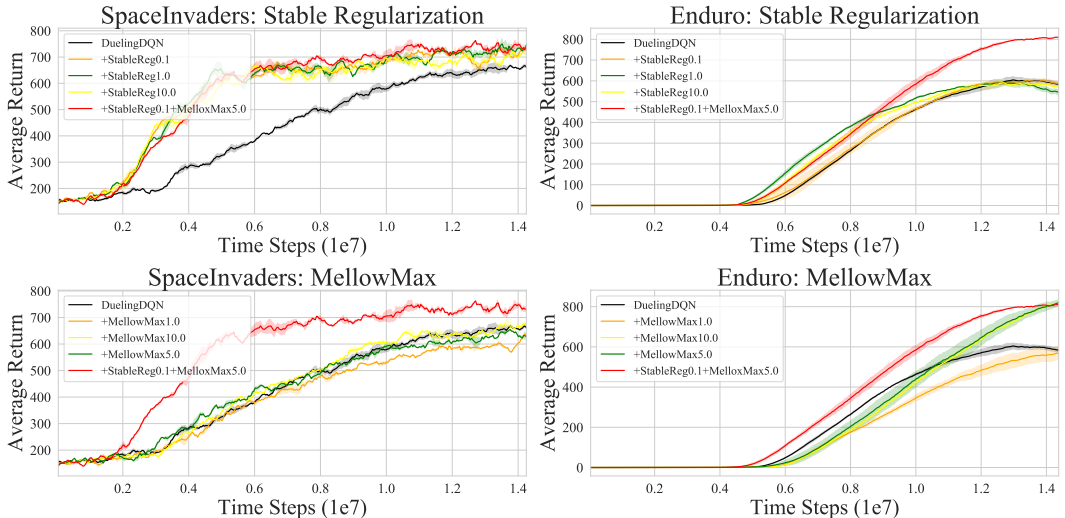


Figure 2: Learning curves of DuelingDQN, +MellowMax, +Stable Regularization (+StableReg), and +Stable Regularization+MellowMax on SpaceInvaders, Enduro games over 3 seeds.

**Impact of Stable Regularization** From diagrams in the first row of Figure 2, we can observe that the benefit margin of stable regularization on Anderson acceleration differs from game to game. Concretely, naively applying stable regularization on DuelingDQN regardless of the MelloxMax to guarantee the faster convergence of Anderson acceleration can still significantly accelerate the convergence in SpaceInvaders. In contrast, the stable regularization is able to boost the sample efficiency mildly on Enduro. For example, when  $\eta = 1.0$  (green line), “+StableReg1.0” is more sample efficient (higher than yellow and orange lines) in the early phase of training. However, the benefit of stable regularization vanishes as the training proceeds, achieving comparable performance with DuelingDQN. Interestingly, if we further add the additional MelloxMax operator (red line), the resulting Stable AA approach can accomplish the improvement of performance to a large margin.

**Impact of MellowMax Operator** As exhibited in diagrams in the last row of Figure 2, the benefit of MelloxMax operator still depends on the game. Particularly, the improvement of MelloxMax operator on SpaceInvaders is negligible, where the lines representing “+MelloxMax” overlap subtly with DuelingDQN. Nevertheless, our Stable AA additionally incorporates the stable regularization, achieving remarkable improvement of sample efficiency. In addition, due to the fact that it is hard to compute the optimal inverse temperature  $\omega$  in MelloxMax, we tune the parameter  $\omega$  and report the corresponding results in Figure 2. It manifests from the diagram on Enduro game that MelloxMax under  $\omega = 5.0, 10.0$  (green and yellow lines) can substantially enhance the final performance. More importantly, under the joint benefits of both the stable regularization and the theoretically-principled MelloxMax operator, our Stable AA DuelingDQN algorithm can simultaneously accelerate the convergence and improve the final performance.

## 5 Discussion and Conclusion

Apart from MelloxMax, other variants of Softmax operator can also be combined with Anderson acceleration, although their theoretical principles have not been studied. For instance, the competitive performance of the Boltzmann Softmax operator suggests that it is still preferable in certain domains,

despite its non-contraction property. We leave the exploration towards more desirable operators as future works. Additionally, the study of our approach on a wider variety of Atari games, and implement on more state-of-the art algorithms are expected in the future.

In this paper, we firstly provide deeper insights into the mechanism of Anderson acceleration on the reinforcement learning setting by connecting damped Anderson acceleration with quasi-Newton method and providing the faster convergence results. These theoretical principles about the faster convergence of Anderson acceleration inspire the leverage of MellowMax operator. Combing with a stable regulation, the resulting Stable AA strategy is applied in DuelingDQN, which has been further demonstrated to significantly accelerate the convergence and enhance the final performance.

## Acknowledgement

We would like to thank the anonymous reviewers for great feedback on the paper. Yingnan Zhao and Ke Sun were supported by the State Scholarship Fund from China Scholarship Council (No:202006120405 and No:202006010082). Dr. Jiang and Dr. Kong were supported by the Natural Sciences and Engineering Research Council of Canada (NSERC). Dr. Kong was also supported by the University of Alberta/Huawei Joint Innovation Collaboration, Huawei Technologies Canada Co., Ltd., and Canada Research Chair in Statistical Learning.

## References

- [1] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [2] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [3] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis, “Mastering the game of go with deep neural networks and tree search,” *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [4] B. Mavrin, S. Zhang, H. Yao, and L. Kong, “Exploration in the face of parametric and intrinsic uncertainties,” in *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, 2019, pp. 2117–2119.
- [5] P. Mirowski, R. Pascanu, F. Viola, H. Soyer, A. J. Ballard, A. Banino, M. Denil, R. Goroshin, L. Sifre, K. Kavukcuoglu *et al.*, “Learning to navigate in complex environments,” *arXiv preprint arXiv:1611.03673*, 2016.
- [6] H. Hasselt, “Double q-learning,” *Advances in Neural Information Processing Systems*, vol. 23, pp. 2613–2621, 2010.
- [7] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas, “Dueling network architectures for deep reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2016, pp. 1995–2003.
- [8] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” *arXiv preprint arXiv:1509.02971*, 2015.
- [9] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 1861–1870.
- [10] M. G. Bellemare, W. Dabney, and R. Munos, “A distributional perspective on reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2017, pp. 449–458.
- [11] B. Mavrin, H. Yao, L. Kong, K. Wu, and Y. Yu, “Distributional reinforcement learning for efficient exploration,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 4424–4434.
- [12] S. Zhang, B. Mavrin, L. Kong, B. Liu, and H. Yao, “Quota: The quantile option architecture for reinforcement learning,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 5797–5804.

- [13] M. L. Puterman, “Markov decision processes,” *Handbooks in Operations Research and Management Science*, vol. 2, pp. 331–434, 1990.
- [14] M. L. Puterman and S. L. Brumelle, “The analytic theory of policy iteration,” *Dynamic Programming and Its Applications*, pp. 91–113, 1978.
- [15] O. Anschel, N. Baram, and N. Shimkin, “Averaged-dqn: Variance reduction and stabilization for deep reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2017, pp. 176–185.
- [16] S. Bubeck, “Convex optimization: Algorithms and complexity,” *arXiv preprint arXiv:1405.4980*, 2014.
- [17] H. F. Walker and P. Ni, “Anderson acceleration for fixed-point iterations,” *SIAM Journal on Numerical Analysis*, vol. 49, no. 4, pp. 1715–1735, 2011.
- [18] A. Fu, J. Zhang, and S. Boyd, “Anderson accelerated douglas–rachford splitting,” *SIAM Journal on Scientific Computing*, vol. 42, no. 6, pp. A3560–A3583, 2020.
- [19] M. Geist and B. Scherrer, “Anderson acceleration for reinforcement learning,” *arXiv preprint arXiv:1809.09501*, 2018.
- [20] Y. Peng, B. Deng, J. Zhang, F. Geng, W. Qin, and L. Liu, “Anderson acceleration for geometry optimization and physics simulation,” *ACM Transactions on Graphics (TOG)*, vol. 37, no. 4, pp. 1–14, 2018.
- [21] H. An, X. Jia, and H. F. Walker, “Anderson acceleration and application to the three-temperature energy equations,” *Journal of Computational Physics*, vol. 347, pp. 1–19, 2017.
- [22] Y. Li, C. Ni, G. Xie, W. Yang, S. Zhou, and Z. Zhang, “Accelerated value iteration via anderson mixing,” 2018.
- [23] W. Shi, S. Song, H. Wu, Y.-C. Hsu, C. Wu, and G. Huang, “Regularized anderson acceleration for off-policy deep reinforcement learning,” *Advances in Neural Information Processing Systems*, 2019.
- [24] H.-r. Fang and Y. Saad, “Two classes of multiseccant methods for nonlinear acceleration,” *Numerical Linear Algebra with Applications*, vol. 16, no. 3, pp. 197–221, 2009.
- [25] C. Evans, S. Pollock, L. G. Rebholz, and M. Xiao, “A proof that anderson acceleration improves the convergence rate in linearly converging fixed-point methods (but not in those converging quadratically),” *SIAM Journal on Numerical Analysis*, vol. 58, no. 1, pp. 788–810, 2020.
- [26] K. Asadi and M. L. Littman, “An alternative softmax operator for reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2017, pp. 243–252.
- [27] M. G. Azar, V. Gómez, and H. J. Kappen, “Dynamic policy programming,” *The Journal of Machine Learning Research*, vol. 13, no. 1, pp. 3207–3245, 2012.
- [28] C. J. C. H. Watkins, “Learning from delayed rewards,” 1989.
- [29] H. Van Hasselt, A. Guez, and D. Silver, “Deep reinforcement learning with double q-learning,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 30, no. 1, 2016.
- [30] L. Pan, Q. Cai, Q. Meng, W. Chen, L. Huang, and T.-Y. Liu, “Reinforcement learning with dynamic boltzmann softmax updates,” *International Joint Conference on Artificial Intelligence (IJCAI)*, 2020.
- [31] Z. Song, R. Parr, and L. Carin, “Revisiting the softmax bellman operator: New benefits and new perspective,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 5916–5925.
- [32] S. Fujimoto, H. Hoof, and D. Meger, “Addressing function approximation error in actor-critic methods,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 1587–1596.
- [33] S. Kim, K. Asadi, M. Littman, and G. Konidaris, “Deepmellow: removing the need for a target network in deep q-learning,” in *Proceedings of the Twenty Eighth International Joint Conference on Artificial Intelligence*, 2019.

## A Proof: Connection with quasi-Newton

*Proof.*

$$\begin{aligned}
Q^{(k+1)} &= (1 - \beta_k) \sum_{l=0}^m \alpha_l^{(k)} Q^{(k-m+l)} + \beta_k \sum_{l=0}^m \alpha_l^{(k)} TQ^{(k-m+l)} \\
&= (1 - \beta_k) \left[ Q^{(k)} - \sum_{i=0}^{m-1} \tau_i \left( Q^{(k-i)} - Q^{(k-i-1)} \right) \right] \\
&\quad + \beta_k \left[ T_{mm}Q^{(k)} - Q^{(k)} + Q^{(k)} - \sum_{i=0}^{m-1} \tau_i \left( TQ^{(k-i)} - T_{mm}Q^{(k-i-1)} \right) \right] \\
&= Q^{(k)} - (1 - \beta_k) \Delta_k \cdot \tau - \beta_k \cdot (\Delta_k + H_k) \tau + \beta_k e_k \\
&= Q^{(k)} + \beta_k e_k - (\Delta_k + \beta_k H_k) \tau \\
&= Q^{(k)} - ((\Delta_k + \beta_k H_k) (H_k^T H_k)^{-1} H_k^T - \beta_k I) e_k \\
&:= Q^{(k)} - G_k e_k
\end{aligned}$$

This formula indicates the term  $G_k = (\Delta_k + \beta_k H_k) (H_k^T H_k)^{-1} H_k^T - \beta_k I$  can be seen as the inverse Jacobian of  $e_k = TQ^{(k)} - Q^{(k)}$ .  $\square$

## B Proof in Convergence results

**Proof about Assumption 1** This proof is to show that MellowMax operator satisfies Assumption 1.

*Proof.* We first show  $T_{mm}$  is twice continuously differentiable. For any vector  $x = (x_1, \dots, x_n)^T$ , we have

$$\frac{\partial mm_\omega(x)}{\partial x_i} = \frac{\exp(\omega x_i)}{\sum_l \exp(\omega x_l)}$$

and

$$\begin{aligned}
\frac{\partial^2 mm_\omega(x)}{\partial x_i^2} &= \frac{\omega \exp(\omega x_i) [\sum_l \exp(\omega x_l)] - (\exp(\omega x_i))^2 \omega}{(\sum_l \exp(\omega x_l))^2} \\
\frac{\partial^2 mm_\omega(x)}{\partial x_i \partial x_j} &= \frac{\partial mm_\omega(x)}{\partial x_i \partial x_j} = \frac{-\omega \exp(\omega(x_i + x_j))}{(\sum_l \exp(\omega x_l))^2},
\end{aligned}$$

which implies that  $T_{mm}$  is twice continuously differentiable due to smoothness  $\exp(\cdot)$  and for any bounded domain in  $\mathbb{R}^n$ , the first and second order derivative exist.

We next show the first and second derivative of  $T_{mm}$  are bounded which follows from  $\|T_{mm}(Q + \Delta) - T_{mm}Q\|_\infty \leq c_1 \|\Delta\|_\infty + c_2 \|\Delta^2\|_\infty + o(\|\Delta^2\|_\infty)$  for any  $\Delta \rightarrow 0$ .

$$\begin{aligned}
\|T_{mm}(Q + \Delta) - T_{mm}Q\|_\infty &= \left\| \frac{\gamma}{\omega} \cdot \mathbf{P} \cdot \log \{ \mathbf{I} \exp[\omega(Q + \Delta)] \} - \frac{\gamma}{\omega} \cdot \mathbf{P} \cdot \log \{ \mathbf{I} \exp(\omega Q) \} \right\|_\infty \\
&= \left\| \frac{\gamma}{\omega} \cdot \mathbf{P} \cdot \log \{ \mathbf{I} \exp(\omega \Delta) \} \right\|_\infty \\
&= \left\| \frac{\gamma}{\omega} \cdot \mathbf{P} \cdot \left[ (\mathbf{I} \exp(\omega \Delta) - \mathbf{I}) - \frac{1}{2} (\mathbf{I} \exp(\omega \Delta) - \mathbf{I})^2 + o(\Delta^2) \right] \right\|_\infty \\
&\leq \left\| \frac{\gamma}{\omega} \cdot \mathbf{P} \cdot [\mathbf{I} \exp(\omega \Delta) - \mathbf{I}] \right\|_\infty + o(\|\Delta^2\|_\infty) \\
&= \left\| \frac{\gamma}{\omega} \cdot \mathbf{P} \cdot \left[ \mathbf{I} \omega \Delta + \frac{1}{2} \mathbf{I} \Delta^2 \omega^2 \right] \right\|_\infty + o(\|\Delta^2\|_\infty) \\
&\leq c_1 \|\Delta\|_\infty + c_2 \|\Delta^2\|_\infty + o(\|\Delta^2\|_\infty), \tag{A1}
\end{aligned}$$

$\square$

where  $\mathbf{P} = [p(s_{i'} | s_i, a_j)]_{1 \leq i, i' \leq |S|, 1 \leq j \leq m}$ .

**Proof about Theorem 1** This proof is to show that we have the results in Theorem 1 under Assumption 1.

*Proof.* Let

$$Q_k^\alpha(s, a) = \sum_{l=0}^m \alpha_l^{(k)} Q^{(k-m+l)}(s, a)$$

$$\widetilde{Q}_k^\alpha(s, a) = \sum_{l=0}^m \alpha_l^{(k)} T_{mm} Q^{(k-m+l)}(s, a)$$

Then

$$Q^{(k+1)}(s, a) = (1 - \beta_k) Q_k^\alpha(s, a) + \beta_k \widetilde{Q}_k^\alpha(s, a).$$

Define  $T'(\cdot; \cdot)$ ,  $T''(\cdot; \cdot, \cdot)$  as linear form with respect to the arguments to the right of semicolon. Let  $\delta_k = Q^{(k)} - Q^{(k-1)}$ ,  $z_k(t) = Q^{(k-1)} + t\delta_k$ ,  $z_{k,t}(u) = z_{k-1}(t) + u(z_k(t) - z_{k-1}(t))$ . Then

$$\begin{aligned} T_{mm}(Q^{(k)}) - T_{mm}(Q^{(k-1)}) &= \int_0^1 T'_{mm}(z_k(t); \delta_k) dt \\ &= \int_0^1 \left\{ T'_{mm}(z_{k+1}(t); \delta_k) + \int_0^1 T''_{mm}(z_{k+1,t}(s); z_k(t) - z_{k+1}(t), \delta_k) ds \right\} dt \\ &= \int_0^1 \int_0^1 \left\{ T'_{mm}(z_{k+1}(t); \delta_k) + T''_{mm}(z_{k+1,t}(s); z_k(t) - z_{k+1}(t), \delta_k) \right\} ds dt. \end{aligned}$$

We note that

$$\begin{aligned} e_k &= T_{mm}(Q^{(k)}) - Q^{(k)} = T_{mm}(Q^{(k)}) - [(1 - \beta_{k-1}) Q_{k-1}^\alpha + \beta_{k-1} \widetilde{Q}_{k-1}^\alpha] \\ &= (1 - \beta_{k-1}) [T_{mm}(Q^{(k)}) - Q_{k-1}^\alpha] + \beta_{k-1} [T_{mm}(Q^{(k)}) - \widetilde{Q}_{k-1}^\alpha] \end{aligned} \quad (\text{A2})$$

For each term on the right hand of formula (A2), we have

$$\begin{aligned} T_{mm} Q^{(k)} - Q_{k-1}^\alpha &= \sum_{i=0}^m \alpha_i^{(k-1)} T_{mm} Q^{(k)} - \sum_{i=0}^m \alpha_i^{(k-1)} Q^{(k-m+i-1)} \\ &= \sum_{i=0}^m \alpha_i^{(k-1)} (T_{mm} Q^{(k)} - Q^{(k-m+i-1)}) \\ &= \sum_{i=0}^m \alpha_i^{(k-1)} (T_{mm} Q^{(k-m+i-1)} - Q^{(k-m+i-1)}) + \sum_{i=0}^m \alpha_i^{(k-1)} (T_{mm} Q^{(k)} - T_{mm} Q^{(k-m+i-1)}) \\ &= e_{k-1}^\alpha + \sum_{i=0}^m \left( \sum_{l=0}^{m-i} \alpha_l^{(k-1)} \right) (T_{mm} Q^{(k-i)} - T_{mm} Q^{(k-i-1)}) \\ &= e_{k-1}^\alpha + \sum_{i=0}^m \tau_i \widetilde{\delta}_{k-i}, \end{aligned}$$

where  $e_k^\alpha = \sum_{i=0}^m \alpha_i^{(k)} (T_{mm} Q^{(k-m+i)} - Q^{(k-m+i)})$ ,  $\tau_i = \sum_{l=0}^{m-i} \alpha_l^{(k-1)}$ ,  $\widetilde{\delta}_{k-i} = T_{mm} Q^{(k-i)} - T_{mm} Q^{(k-i-1)}$ . Moreover,

$$\begin{aligned} T_{mm} Q^{(k)} - \widetilde{Q}_{k-1}^\alpha &= T_{mm} Q^{(k)} - \sum_{i=0}^m \alpha_i^{(k-1)} T_{mm} Q^{(k-i-1)} \\ &= \sum_{i=0}^m \alpha_i^{(k-1)} (T_{mm} Q^{(k)} - T_{mm} Q^{(k-i-1)}) \\ &= \sum_{i=0}^m \tau_i \widetilde{\delta}_{k-i}. \end{aligned}$$

Therefore, formula (A2) can be rewritten as

$$\begin{aligned}
e_k &= (1 - \beta_{k-1})(e_{k-1}^\alpha + \sum_{i=0}^m \tau_i \widetilde{\delta_{k-i}}) + \beta_{k-1} \sum_{i=0}^m \tau_i \widetilde{\delta_{k-i}} \\
&= (1 - \beta_{k-1})e_{k-1}^\alpha + \sum_{i=0}^m \tau_i \widetilde{\delta_{k-i}} \\
&= (1 - \beta_{k-1})e_{k-1}^\alpha + \sum_{i=0}^m \tau_i \int_0^1 T'_{mm}(z_{k-i}(t); \delta_{k-i}) dt \\
&= (1 - \beta_{k-1})e_{k-1}^\alpha + \sum_{i=1}^m \tau_i \left\{ \int_0^1 T'_{mm}(z_k(t); \delta_{k-i}) dt \right. \\
&\quad \left. + \sum_{l=k-i}^{k-1} \int_0^1 T'_{mm}(z_l(t); \delta_{k-i}) - T'_{mm}(z_{l+1}(t); \delta_{k-i}) dt \right\} + \int_0^1 T'_{mm}(z_k(t); \delta_k) dt \\
&= (1 - \beta_{k-1})e_{k-1}^\alpha + \int_0^1 T'_{mm}(z_k(t); \sum_{i=0}^m \tau_i \delta_{k-i}) dt \\
&\quad + \sum_{i=1}^m \tau_i \sum_{l=k-i}^{k-1} \int_0^1 \int_0^1 T''_{mm}(z_{l+1,t}(s); z_l(t) - z_{l+1}(t), \delta_{k-i}) ds dt \\
&= (1 - \beta_{k-1})e_{k-1}^\alpha + \int_0^1 T'_{mm}(z_k(t); \sum_{i=0}^m \tau_i \delta_{k-i}) dt \\
&\quad + \sum_{i=1}^m \int_0^1 \int_0^1 \sum_{l=k-i}^{k-1} T''_{mm}(z_{l+1,t}(s); z_l(t) - z_{l+1}(t), \tau_i \delta_{k-i}) ds dt.
\end{aligned}$$

For the term  $\sum_{i=0}^m \tau_i \delta_{k-i}$ , it can be rewritten as

$$\begin{aligned}
\sum_{i=0}^m \tau_i \delta_{k-i} &= \delta_k + \sum_{i=1}^m \tau_i \delta_{k-i} \\
&= Q^{(k)} - Q^{(k-1)} + \tau_1 Q^{(k-1)} - \sum_{i=0}^{m-1} \alpha_i Q^{(k-m+i-1)} \\
&= Q^{(k)} - \alpha_m^{(k-1)} Q^{(k-1)} - \sum_{i=1}^{m-1} \alpha_i^{(k-1)} Q^{(k-m+i-1)} \\
&= Q^{(k)} - Q_{k-1}^\alpha \\
&= \beta_{k-1} (\widetilde{Q_{k-1}^\alpha} - Q_{k-1}^\alpha) = \beta_{k-1} e_{k-1}^\alpha,
\end{aligned}$$

where the second and third equality hold using the formula  $\tau_i - \tau_{i+1} = \alpha_{m-i}^{(k-1)}$ ,  $\tau_1 = 1 - \alpha_m^{(k-1)}$ . Then, we obtain

$$\begin{aligned}
e_k &= \int_0^1 (1 - \beta_{k-1})e_{k-1}^\alpha + \beta_{k-1} T'_{mm}(z_k(t); e_{k-1}^\alpha) dt \\
&\quad + \sum_{i=1}^m \int_0^1 \int_0^1 \sum_{l=k-i}^{k-1} T''_{mm}(z_{l+1,t}(s); z_l(t) - z_{l+1}(t), \tau_i \delta_{k-i}) ds dt. \quad (\text{A3})
\end{aligned}$$

Formula (A1) and (A3) together imply that

$$\begin{aligned}
\|e_k\|_\infty &\leq (1 - \beta_{k-1}) \|e_{k-1}^\alpha\|_\infty + \beta_{k-1} \cdot c_1 \cdot \|e_{k-1}^\alpha\|_\infty + \sum_{i=1}^m \sum_{l=k-i}^{k-1} c_2 \cdot (\|\delta_l\|_\infty + \|\delta_{l+1}\|_\infty) |\tau_i| \|\delta_{k-i}\|_\infty \\
&= \theta_k \left\{ ((1 - \beta_{k-1}) + c_1 \beta_{k-1}) \|e_{k-1}^\alpha\|_\infty \right\} + c_2 \cdot \sum_{i=2}^m \left( \|\delta_k\|_\infty + \|\delta_{k-i}\|_\infty + 2 \sum_{l=1}^{i-1} \|\delta_{k-i}\|_\infty \right) |\tau_i| \|\delta_{k-i}\|_\infty \\
&\quad + c_2 \cdot (\|\delta_k\|_\infty + \|\delta_{k-1}\|_\infty) |\tau_1| \|\delta_{k-1}\|_\infty. \quad (\text{A4})
\end{aligned}$$

□

**Proof about Assumption 2** This proof is to show that MellowMax operator satisfies Assumption 2 (non-expansive operator). Similar result is also given in [26, 33].

*Proof.* Let  $|\mathcal{S}| = n_1, \mathcal{A} = n_2$ . Note that

$$T_{mm}Q = R + \gamma \cdot \mathbf{P} \cdot mm_\omega(Q)$$

where  $mm_\omega(Q) = \frac{1}{\omega} \log\{\frac{1}{n_2} \cdot \mathbf{I} \cdot \exp(\omega Q)\}$ ,  $\mathbf{I} = \mathbf{I}_{n_1 \times n_1} \otimes \mathbf{1}_{n_2 \times 1}^T$ .

$$\begin{aligned} \|T_{mm}Q - T_{mm}Q'\|_\infty &\leq \gamma \|\mathbf{P}\|_\infty \|mm_\omega(Q) - mm_\omega(Q')\|_\infty \\ &\leq \gamma \|mm_\omega(Q) - mm_\omega(Q')\|_\infty \\ &\leq \gamma \|Q - Q'\|_\infty \end{aligned} \quad (\text{A5})$$

□

**Proof about Theorem 2** We analyze a bound for  $\delta_j$  in terms of  $e_j$  in the following part. Based on formula (A5), we have

$$\begin{aligned} (1 - \gamma) \|\delta_k\|_\infty &= \|\delta_k\|_\infty - \gamma \|\delta_k\|_\infty \\ &\leq \|\delta_k\|_\infty - \|T_{mm}Q^{(k)} - T_{mm}Q^{(k-1)}\|_\infty \\ &\leq \|Q^{(k)} - Q^{(k-1)} - T_{mm}Q^{(k)} + T_{mm}Q^{(k-1)}\|_\infty \\ &= \|e_k - e_{k-1}\|_\infty. \end{aligned} \quad (\text{A6})$$

Let  $E_k = (e_{k-m}, \dots, e_k)$ . The optimization problem

$$\alpha^k = \operatorname{argmin}_{\alpha \in \mathbb{R}^{m+1}} \|E_k \alpha\|_2^2 \quad \text{s.t.} \quad \sum_{i=0}^m \alpha_i = 1$$

is equivalent to the unconstrained form

$$\min_{\eta \in \mathbb{R}^m} \|e_{k-m} + \sum_{i=1}^m \eta_i (e_{k-m+i} - e_{k-h+i-1})\|^2, \quad \eta_i = \sum_{l=i}^m \alpha_l^{(k)} \quad (\text{A7})$$

$$\min_{\tilde{\tau} \in \mathbb{R}^m} \left\| e_k - \sum_{i=0}^{m-1} \tilde{\tau}_i (e_{k-i} - e_{k-i-1}) \right\|^2, \quad \tilde{\tau}_i = \sum_{l=0}^{m-i-1} \alpha_l^{(k)} \quad (\text{A8})$$

Seeking the critical point for  $\eta_m$  in (A7) yields that

$$\langle e_{k-m}, e_k - e_{k-1} \rangle + \sum_{i=1}^m \eta_i \langle e_{k-m+i} - e_{k-m+i-1}, e_k - e_{k-1} \rangle = 0.$$

This implies that

$$\begin{aligned} \eta_m \|e_k - e_{k-1}\|^2 &= -\langle e_{k-m}, e_k - e_{k-1} \rangle - \sum_{i=1}^{m-1} \eta_i \langle e_k - e_{k-1}, e_{k-m+i} - e_{k-m+i-1} \rangle \\ &= -\eta_{m-1} \langle e_k - e_{k-1}, e_{k-1} \rangle - \left\langle e_k - e_{k-1}, \sum_{i=0}^{m-2} \alpha_i e_{k-m+i} \right\rangle. \end{aligned}$$

Applying Cauchy-Schwarz inequality and triangle inequalities yields

$$\left| \alpha_m^{(k)} \right| \|e_k - e_{k-1}\| \leq |\eta_{m-1}| \|e_{k-1}\| + \sum_{i=0}^{m-2} \alpha_i^{(k)} \|e_{k-m+i}\|.$$

Based on the inequality  $\|\cdot\|_\infty \leq \|\cdot\|_2$  over  $\mathbb{R}^n$  and formula (A6), it follows

$$\left| \alpha_m^{(k)} \right| \|\delta_k\|_\infty \leq \frac{1}{1 - \gamma} \left\{ |\eta_{m-1}| \|e_{k-1}\| + \sum_{i=0}^{m-2} \alpha_i^{(k)} \|e_{k-m+i}\| \right\}. \quad (\text{A9})$$

Seeking the critical point with respect to  $\tilde{\tau}_p$ , ( $p = 1, \dots, m-1$ ) in (A8) yields

$$\left\langle e_k - \sum_{i=0}^{m-1} \tilde{\tau}_i (e_{k-i} - e_{k-i-1}), e_{k-p} - e_{k-p-1} \right\rangle = 0$$

which implies

$$\begin{aligned} \tilde{\tau}_p \|e_{k-p} - e_{k-p-1}\|^2 &= \langle e_{k-p} - e_{k-p-1}, \tilde{\tau}_{p-1} e_{k-p} \rangle - \langle e_{k-p} - e_{k-p-1}, \tilde{\tau}_{p+1} e_{k-p-1} \rangle \\ &\quad + \left\langle e_{k-p} - e_{k-p-1}, \sum_{j=0}^{m-p-2} \alpha_j e_{k-m+j} \right\rangle + \left\langle e_{k-p} - e_{k-p-1}, \sum_{j=m-p+1}^m \alpha_j e_{k-m+j} \right\rangle. \end{aligned}$$

Then

$$|\tilde{\tau}_p| \|\delta_{k-p}\|_\infty \leq \frac{1}{1-\gamma} \left\{ |\tilde{\tau}_{p-1}| \|e_{k-p}\| + |\tilde{\tau}_{p+1}| \|e_{k-p-1}\| + \sum_{j=0}^{m-p-2} |\alpha_j| \|e_{k-m+j}\| + \sum_{j=m-p+1}^m |\alpha_j| \|e_{k-m+j}\| \right\} \quad (\text{A10})$$

Combing (A4), (A9) and (A10), we establish

$$\begin{aligned} \|e_k\|_\infty &\leq \theta_k \left\{ ((1 - \beta_{k-1}) + c_1 \beta_{k-1}) \|e_{k-1}\|_\infty \right\} + \text{Constant} \cdot \left\{ \sum_{i=2}^m \|\delta_{k-i}\|_\infty^2 + \|\delta_{k-1}\|_\infty^2 \right\} \\ &= \theta_k \left\{ ((1 - \beta_{k-1}) + c_1 \beta_{k-1}) \|e_{k-1}\|_\infty \right\} + O \left( \sum_{i=1}^m \|e_{k-i}\|_\infty^2 \right). \end{aligned}$$

## C Proof: Stable regularization

We firstly prove the stability of the derived regularization in Theorem 3.

*Proof.* It is easy to prove that  $\|\tilde{G}_k^{-1} G_k^{-1}\|_2 \leq 1$  as long as we directly remove the regularization term to induce the inequality. Then, we have

$$\begin{aligned} \|\tilde{G}_k\|_2 &\leq \left| -\beta_k + \frac{\|\Delta_k + \beta_k H_k\|_2 \|H_k\|_2}{\eta (\|\Delta_k\|_F^2 + \|H_k\|_F^2)} \right| \\ &\leq \left| -\beta_k + \frac{\|\Delta_k\|_2 \|H_k\|_2 + \|H_k\|_2^2}{\eta (\|\Delta_k\|_F^2 + \|H_k\|_F^2)} \right| \\ &\leq \left| -\beta_k + \frac{\|\Delta_k\|_F \|H_k\|_F + \|H_k\|_F^2}{\eta (\|\Delta_k\|_F^2 + \|H_k\|_F^2)} \right| \\ &\leq \left| \frac{2}{\eta} - \beta_k \right|. \end{aligned}$$

□

This indicates that the  $\ell_2$  norm of updating matrix  $\tilde{G}_k$  is upper bounded, which can guarantee the stability. Then we provide the proof of Proposition 2.

*Proof.* Firstly, we denote the structure matrix A as follows:

$$A = \begin{pmatrix} 0 & 0 & 0 & 0 & \cdots & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & \cdots & 0 & 1 & -1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 1 & -1 & \cdots & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & \cdots & 0 & 0 & 0 \\ 1 & -1 & 0 & 0 & \cdots & 0 & 0 & 0 \end{pmatrix}_{(m+1) \times (m+1)}$$

Note that  $\alpha_{\text{reg}}^{(k)} = A \cdot \tilde{\tau}_{\text{reg}}$ , where

$$\tilde{\tau}_{\text{reg}} = \begin{pmatrix} 1 \\ \tau_{\text{reg}} \end{pmatrix} \in \mathbb{R}^{m+1}.$$

We first bound  $\alpha_{\text{reg}}^{(k)}$ ,



$$\begin{aligned}
\|\alpha_{\text{reg}}^{(k)}\|_2^2 &\leq \|A\|_2^2 \cdot \|\tilde{\tau}_{\text{reg}}\|_2^2 \leq 4 \cdot (1 + \|\tau_{\text{reg}}\|_2^2) \\
&\leq 4 \left[ 1 + \left\| \left( H_k^T H_k + \eta (\|\Delta_k\|_F^2 + \|H_k\|_F^2) I \right)^{-1} \right\|^2 \cdot \|H_k^T e_k\|^2 \right] \\
&\leq 4 \left( 1 + \frac{\|H_k^T e_k\|^2}{\eta^2 (\|\Delta_k\|_F^2 + \|H_k\|_F^2)} \right) \\
&\leq 4 \left( 1 + \frac{\|e_k\|^2}{\eta^2} \right).
\end{aligned}$$

We next analyze  $\alpha_{\text{reg}}^{(k)} - \alpha_{\text{non}}^{(k)}$ . Since

$$\begin{aligned}
H_k^T e_k - H_k^T H_k \tau_{\text{non}} &= 0, \\
H_k^T e_k - \left[ H_k^T H_k + \eta (\|\Delta_k\|_F^2 + \|H_k\|_F^2) I \right] \tau_{\text{reg}} &= 0
\end{aligned}$$

Then  $\tau_{\text{reg}} - \tau_{\text{non}} = \left[ H_k^T H_k + \eta (\|\Delta_k\|_F^2 + \|H_k\|_F^2) I \right]^{-1} \left[ \eta (\|\Delta_k\|_F^2 + \|H_k\|_F^2) I \right] \tau_{\text{non}}$  which implies

$$\|\tau_{\text{reg}} - \tau_{\text{non}}\|_2 \leq \frac{(\eta (\|\Delta_k\|_F^2 + \|H_k\|_F^2)) \|I\|_2}{\eta (\|\Delta_k\|_F^2 + \|H_k\|_F^2)} \|\tau_{\text{non}}\|_2 = \|\tau_{\text{non}}\|_2. \quad (\text{A11})$$

Let  $\tilde{\tau}_{\text{reg}} = (1, \tau_{\text{reg}}^T)^T$ ,  $\tilde{\tau}_{\text{non}} = (1, \tau_{\text{non}}^T)^T$ . Then  $\tilde{\tau}_{\text{non}} = A^{-1} \alpha_{\text{non}}^{(k)}$ , and  $\|\tau_{\text{non}}\|_2^2 = \|A^{-1} \alpha_{\text{non}}^{(k)}\|_2^2 - 1$ . Based on (A11), we can establish

$$\begin{aligned}
\|\alpha_{\text{reg}}^{(k)} - \alpha_{\text{non}}^{(k)}\|_2^2 &\leq \|A\|_2^2 \|\tilde{\tau}_{\text{reg}} - \tilde{\tau}_{\text{non}}\|_2^2 = \|A\|_2^2 \|\tau_{\text{reg}} - \tau_{\text{non}}\|_2^2 \\
&\leq \|A\|_2^2 \left( \|A^{-1} \alpha_{\text{non}}^{(k)}\|_2^2 - 1 \right) \\
&\leq \|A\|_2^2 \cdot \|A^{-1}\|_2^2 \cdot \|\alpha_{\text{non}}^{(k)}\|_2^2 - \|A\|_2^2 \\
&\leq (\text{cond}_2(A))^2 \cdot \|\alpha_{\text{non}}^{(k)}\|_2^2 - \frac{2m+1}{m+1}.
\end{aligned}$$

□

## D Results on Other games

We provide results of our algorithms on other 12 Atari games. Our results in Figure 3,4,5 and 6 show that our Stable AA DuelingDQN consistently outperforms both DuelingDQN and DuelingDQN-RAA significantly.

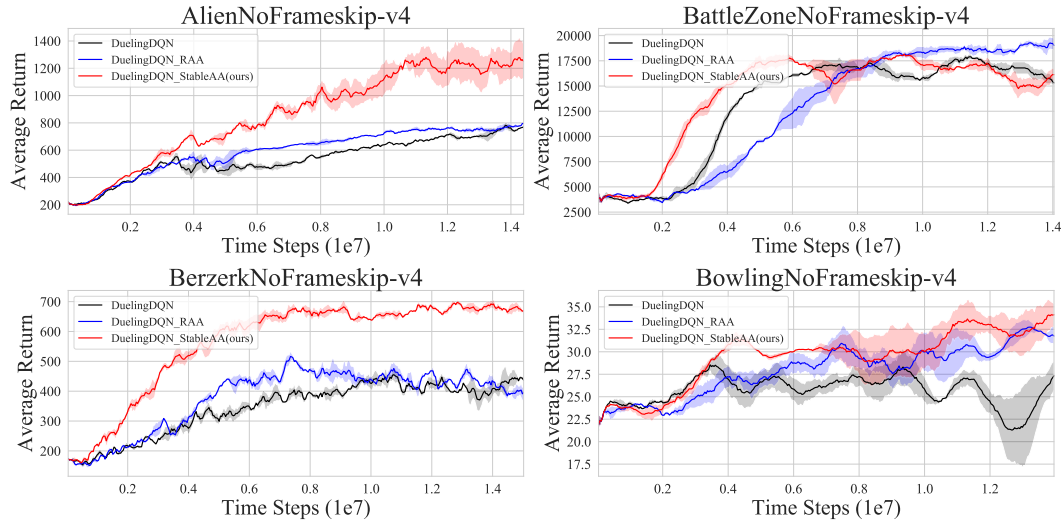


Figure 3: Learning curves of DuelingDQN, DuelingDQN-RAA, DuelingDQN-Stable AA (ours) on Alien, BattleZone, Berzerk and Bowling games over 3 seeds.

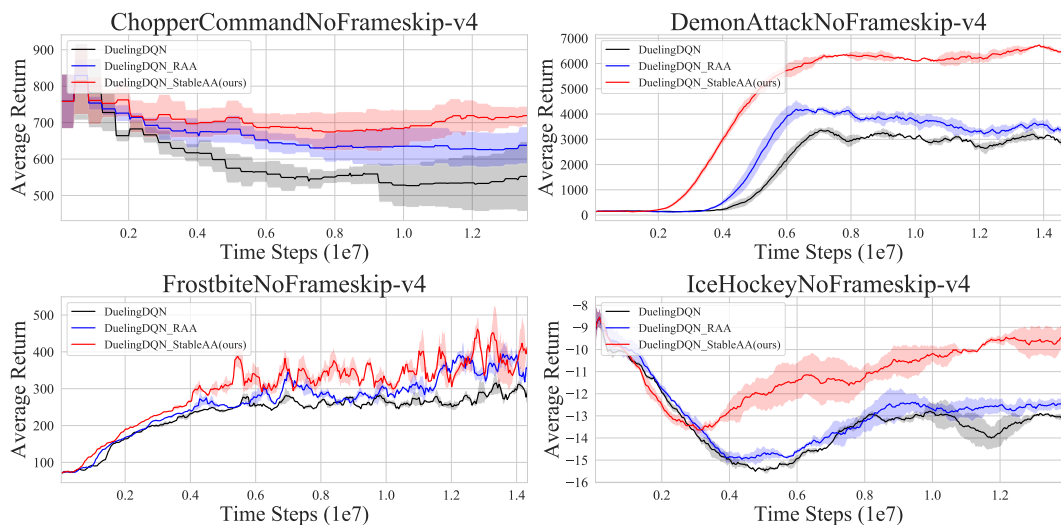


Figure 4: Learning curves of DuelingDQN, DuelingDQN-RAA, DuelingDQN-Stable AA (ours) on ChopperCommand, DemonAttack, Frostbite and IceHockey games over 3 seeds.

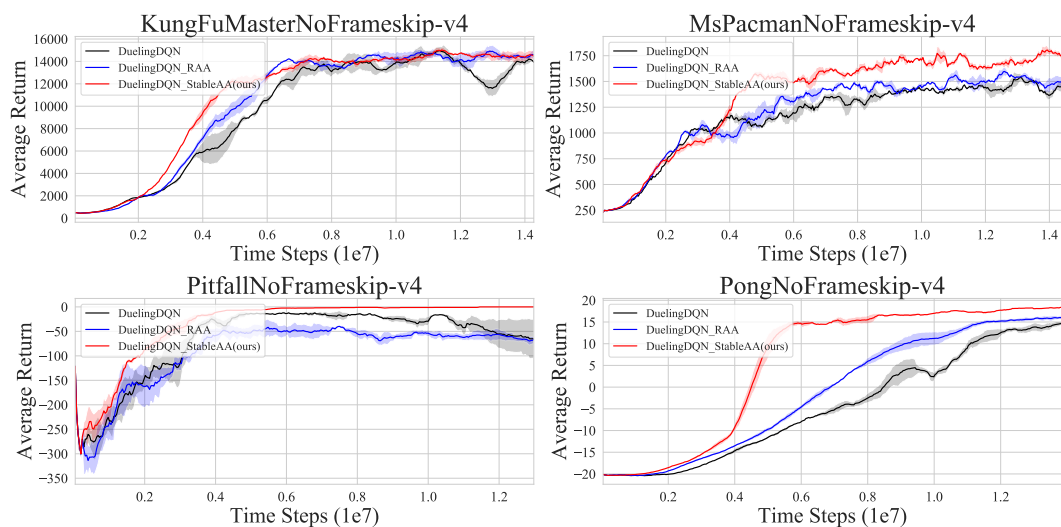


Figure 5: Learning curves of DuelingDQN, DuelingDQN-RAA, DuelingDQN-Stable AA (ours) on KungFu, MsPacman, Pitfall and Pong games over 3 seeds.

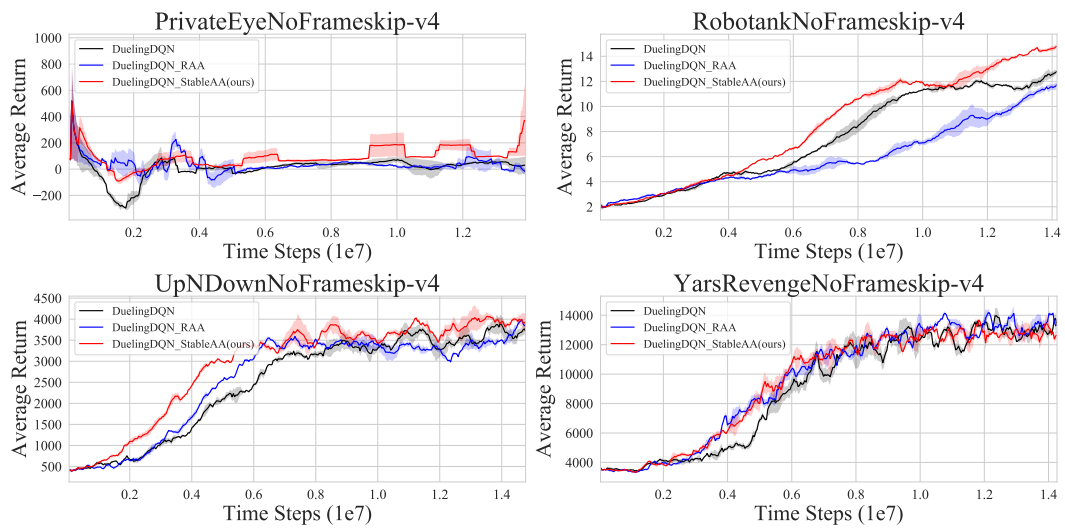


Figure 6: Learning curves of DuelingDQN, DuelingDQN-RAA, DuelingDQN-Stable AA (ours) on PrivateEye, Robotank, UpNDown and YarsRevenge games over 3 seeds.