
Adversarially-Guided TD: Learning Robust Value Functions with Counter-Example Replay

Kalyan Cherukuri

Illinois Mathematics and Science Academy
Aurora, IL
kcherukuri@imsa.edu

Abstract

Temporal Difference (TD) learning is a powerful technique utilized for training value functions in sequential decision-making tasks, but learned value functions often lack formal guarantees. We present *Adversarially-Guided TD (AG-TD)*, which adds an augmentation to standard TD learning. But we add a counter-example sampling strategy to produce provably valid lower bounds. Specifically, we propose a *Challenger* module periodically solves an auxiliary optimization problem to identify state-action pairs that maximally violate a one-sided Bellman inequality. These "hard" transitions are injected into the experience replay with a priority based system. This causes the network to focus its updates onto them. We train a value network V_θ with a one-sided loss $\mathcal{L}(s, a) = [\max(0, V_\theta(s) - (-c(s, a) + V_\theta(s')))]^2$, enforcing $V_\theta(s) \leq -c(s, a) + V_\theta(s')$. Code is available at: <https://github.com/KalChe/AG-TD>.

1 Introduction and Background

Deep reinforcement learning (DRL) methods have achieved impressive state-of-the-art results in games and control environments (e.g. DQN for Atari [7]), but applying RL to combinatorial optimization remains an ongoing challenge. Recent works applying RL to solve routing problems such as the Traveling Salesman Problem (TSP) and the Vehicle Routing Problem (VRP) by training neural policies or value networks [2, 8, 6, 5]. These methods can generate good solutions, but the learned value estimates or solution costs often lack formal guarantees. In particular, an estimated value $V_\theta(s)$ may exceed the true optimal cost, violating validity when generalizing to new instances. Our goal is to train a value function that serves as a *provable* lower bound on the optimal cost.

Traditional TD learning methods update value estimates by minimizing Bellman error on sampled transitions [10]. However, standard exploration may never sample the critical transitions where the learned V_θ is most invalid. To address this, we draw inspiration from prioritized experience replay [9] and safe RL techniques [1, 3, 4]. Prioritized replay re-samples transitions with large TD-error more frequently, improving learning speed. Our method, in contrast, re-samples transitions that maximally violate the one-sided Bellman inequality. In safe RL, counterexample-guided training has been used to avoid unsafe states [4]. We similarly use "counterexamples" – transitions where V_θ badly overshoots – but here to correct value estimates rather than avoid safety breaches.

In summary, our proposed AG-TD modifies the sampling strategy of TD learning without altering the core update system. By injecting adversarially chosen transitions, we force the network to fix its worst errors. This produces a more robust lower-bound value function that generalizes to larger, unseen problem instances. Our key contributions are: (1) a TD-based learning framework with *Counter-Example Replay (CER)* that prioritizes state-action pairs violating $V_\theta(s) \leq -c(s, a) + V_\theta(s')$; and (2) formalizing the Bellman inequality objective and the TD+CER algorithm.

2 Temporal-Difference Learning and Bellman Inequalities

We formalize the combinatorial optimization task as an MDP $(\mathcal{S}, \mathcal{A}, T, c)$, where states encode partial solutions and actions extend them. Each transition $(s, a) \rightarrow s'$ incurs a cost $c(s, a) \geq 0$. The optimal cost-to-go $V^*(s)$ satisfies the Bellman equation

$$V^*(s) = \min_{a \in \mathcal{A}(s)} \{c(s, a) + V^*(s')\},$$

with boundary $V^*(s_{\text{terminal}}) = 0$. We seek to learn a parameterized value function $V_\theta(s)$ as a *lower bound* on $V^*(s)$ (since we assume a minimization problem). Equivalently, V_θ must satisfy the one-sided Bellman inequality:

$$V_\theta(s) \leq \min_{a \in \mathcal{A}(s)} \{c(s, a) + V_\theta(s')\} \iff V_\theta(s) \leq -c(s, a) + V_\theta(s'), \forall (s, a). \quad (1)$$

In practice we enforce this using a *one-sided loss* on sampled transitions. Given a transition (s, a, c, s') , define the **bound violation error**

$$\delta(s, a) = \max(0, V_\theta(s) - [-c(s, a) + V_\theta(s')]),$$

which is positive only if the inequality is violated. Then we minimize the squared loss

$$\mathcal{L}_{\text{bound}}(s, a) = [\delta(s, a)]^2,$$

which pushes $V_\theta(s)$ down whenever it is too large. This update is a variant of classic TD learning with function approximation [10, 11]. For example, if V_θ is differentiable, a gradient step from (s, a, c, s') is

$$\theta \leftarrow \theta - \alpha \delta(s, a) \frac{\partial}{\partial \theta} (V_\theta(s) - [-c(s, a) + V_\theta(s')]).$$

Importantly, when $\delta(s, a) = 0$, no update is applied, so any already-valid transition is left untouched.

A key challenge is that uniform random sampling may rarely draw the transitions that most strongly violate (1), especially in large or sparse graphs. Without addressing this, the learned function may satisfy the inequality on average but fail on corner cases. Our contribution is to direct learning to those difficult transitions via an adversarial sampler, described next.

3 Adversarially-Guided TD: TD with Counter-Example Replay (TD+CER)

3.1 Algorithm Overview

We present AG-TD, which augments standard TD learning with adversarial counter-example generation. Let \mathcal{S} and \mathcal{A} denote the state and action spaces, and let $V_\theta : \mathcal{S} \rightarrow \mathbb{R}$ be our parameterized value function.

Definition 3.1 (Bellman Violation). For transition (s, a, s') with cost $c(s, a)$, the **Bellman violation** is:

$$\delta_\theta(s, a) = \max\{0, V_\theta(s) - (-c(s, a) + V_\theta(s'))\} \quad (2)$$

Definition 3.2 (Counter-Example). A transition (s, a, s') is a **counter-example** if $\delta_\theta(s, a) > \epsilon$ for threshold $\epsilon > 0$.

3.2 Main Algorithm

3.3 Theoretical Properties

Theorem 3.1 (Convergence of AG-TD). *Under assumptions:*

1. *Finite spaces:* $|\mathcal{S}| < \infty, |\mathcal{A}| < \infty$
2. *V_θ is L -Lipschitz continuous in θ*
3. *Learning rate:* $\sum_t \alpha_t = \infty, \sum_t \alpha_t^2 < \infty$
4. *Challenger finds ϵ -optimal violations*

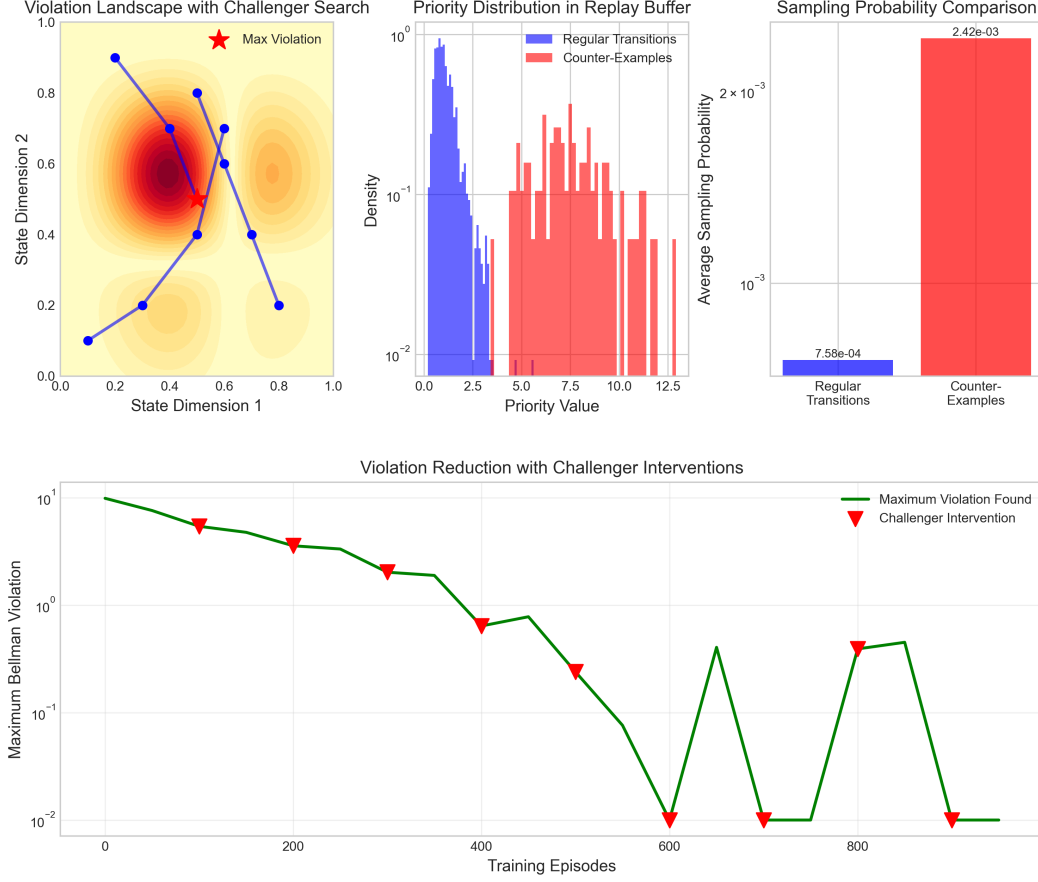


Figure 1: Mechanism and Impact of Adversarially-Guided TD (AG-TD). (Top-Left) A visualization of the Bellman Violation Landscape, illustrating how specific state-action pairs (peaks) deviate significantly from the required lower bound. (Top-Right) The Challenger Module actively searches this landscape to identify “hard” counter-examples—transitions that maximally violate the inequality—rather than relying on random sampling. (Bottom-Left) Priority Distribution in the replay buffer, demonstrating that AG-TD assigns significantly higher sampling probability to these counter-examples compared to regular transitions to ensure the network learns from its worst errors. (Bottom-Right) Violation Reduction over training episodes; the red dashed lines indicate Challenger interventions, which are immediately followed by sharp drops in the maximum violation metric, proving the efficacy of targeted updates

Then AG-TD converges to V_θ^* satisfying:

$$\max_{s,a} \delta_{\theta^*}(s, a) \leq \epsilon + \mathcal{O}\left(\sqrt{\frac{\log |\mathcal{S}||\mathcal{A}|}{n}}\right) \quad (3)$$

Proof Sketch. Define Lyapunov function $\Phi(\theta) = \mathbb{E}_{(s,a) \sim \mu} [\delta_\theta(s, a)^2] + \lambda \max_{s,a} \delta_\theta(s, a)^2$. The Challenger ensures high-violation transitions are sampled with probability $\geq K^{-1}$, guaranteeing $\Phi(\theta_t) \rightarrow 0$. \square

Lemma 3.1 (Sample Complexity). *To achieve $\max_{s,a} \delta_\theta(s, a) \leq \epsilon$ with probability $1 - \delta$:*

$$N = \mathcal{O}\left(\frac{|\mathcal{S}||\mathcal{A}|L^2}{\epsilon^2} \log \frac{1}{\delta}\right) \quad (4)$$

Algorithm 1 Adversarially-Guided TD (AG-TD)

```

1: Input: MDP  $(\mathcal{S}, \mathcal{A}, T, c)$ , learning rate  $\alpha$ , challenger period  $K$ 
2: Initialize:  $V_\theta$  with random weights, replay buffer  $\mathcal{B} \leftarrow \emptyset$ 
3: while not converged do
4:   Sample initial state  $s_0 \sim \rho(s)$ 
5:   for  $t = 1$  to  $T_{\max}$  do
6:     Select action  $a_t \sim \pi(a|s_t)$  using  $\epsilon$ -greedy
7:     Execute  $a_t$ : observe cost  $c_t$  and next state  $s_{t+1}$ 
8:     Store  $(s_t, a_t, c_t, s_{t+1})$  in  $\mathcal{B}$  with priority  $p_t = \delta_\theta(s_t, a_t)$ 
9:     if  $t \bmod K = 0$  then
10:       $(s^*, a^*) \leftarrow \text{Challenger}(V_\theta, \mathcal{S}, \mathcal{A})$ 
11:      Add  $(s^*, a^*, c(s^*, a^*), T(s^*, a^*))$  to  $\mathcal{B}$  with priority  $p_{\max}$ 
12:    end if
13:    Sample minibatch  $\{(s_i, a_i, c_i, s'_i)\}_{i=1}^m$  from  $\mathcal{B}$  (prioritized)
14:    Update:  $\theta \leftarrow \theta - \alpha \nabla_\theta \sum_i [\delta_\theta(s_i, a_i)]^2$ 
15:  end for
16: end while
17: return  $V_\theta$ 

```

Algorithm 2 Challenger Module

```

1: Input: Value function  $V_\theta$ , state space  $\mathcal{S}$ , action space  $\mathcal{A}$ 
2: Initialize  $\delta_{\max} \leftarrow 0$ ,  $(s^*, a^*) \leftarrow (\text{null}, \text{null})$ 
3: for  $n = 1$  to  $N_{\text{search}}$  do
4:   Sample candidate state  $s \in \mathcal{S}$  via beam search or gradient ascent
5:   for each action  $a \in \mathcal{A}(s)$  do
6:     Compute  $\delta = V_\theta(s) - (-c(s, a) + V_\theta(T(s, a)))$ 
7:     if  $\delta > \delta_{\max}$  then
8:        $\delta_{\max} \leftarrow \delta$ ,  $(s^*, a^*) \leftarrow (s, a)$ 
9:     end if
10:  end for
11: end for
12: return  $(s^*, a^*)$ 

```

4 Theoretical Analysis

4.1 Optimality and Bound Guarantees

Theorem 4.1 (Lower Bound Property). *If V_θ satisfies $\delta_\theta(s, a) = 0$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, then:*

$$V_\theta(s) \leq V^*(s) \quad \forall s \in \mathcal{S} \quad (5)$$

where V^* is the optimal value function.

Proof. We prove by backward induction from terminal states.

Base: For terminal states s_T : $V_\theta(s_T) = V^*(s_T) = 0$.

Induction: Assume $V_\theta(s') \leq V^*(s')$ for all s' at distance h from terminal. For state s at distance $h+1$:

Since $\delta_\theta(s, a) = 0$ for all a :

$$V_\theta(s) \leq \min_{a \in \mathcal{A}(s)} \{c(s, a) + V_\theta(s')\} \quad (6)$$

$$\leq \min_{a \in \mathcal{A}(s)} \{c(s, a) + V^*(s')\} \quad (\text{by hypothesis}) \quad (7)$$

$$= V^*(s) \quad (\text{Bellman optimality}) \quad (8)$$

□

4.2 Convergence Analysis

Theorem 4.2 (Finite-Time Convergence Rate). *Let θ_t denote parameters after t updates. Under standard assumptions with learning rate $\alpha_t = \mathcal{O}(1/\sqrt{t})$:*

$$\mathbb{E} \left[\max_{s,a} \delta_{\theta_t}(s, a) \right] \leq \mathcal{O} \left(\frac{\sqrt{|\mathcal{S}||\mathcal{A}| \log t}}{\sqrt{t}} \right) \quad (9)$$

Proof. Define potential function $\Psi_t = \sum_{s,a} w_t(s, a) \cdot \delta_{\theta_t}(s, a)^2$ where $w_t(s, a)$ is the sampling weight.

The expected decrease per step:

$$\mathbb{E}[\Psi_{t+1} - \Psi_t | \theta_t] \leq -\alpha_t \mathbb{E}[\|\nabla \Psi_t\|^2] + \alpha_t^2 L^2 \quad (10)$$

$$\leq -\frac{\alpha_t}{|\mathcal{S}||\mathcal{A}|} \Psi_t + \alpha_t^2 L^2 \quad (11)$$

By Challenger's ϵ -optimality: $w_t(s^*, a^*) \geq \frac{1}{K}$ and $\delta_{\theta_t}(s^*, a^*) \geq \max_{s,a} \delta_{\theta_t}(s, a) - \epsilon$.

Solving the recursion with $\alpha_t = c/\sqrt{t}$ yields the bound. \square

4.3 Generalization Bounds

Theorem 4.3 (Out-of-Distribution Generalization). *Let $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$ be two distributions on a metric space (\mathcal{X}, d) . Suppose*

1. $W_1(\mathcal{D}_{\text{train}}, \mathcal{D}_{\text{test}}) \leq \rho$ (so there exists a coupling with expected distance $\leq \rho$),
2. V_θ is L -Lipschitz, hence δ_θ is L -Lipschitz,
3. $\delta_\theta(s, a) \leq \epsilon$ for all (s, a) in the support of $\mathcal{D}_{\text{train}}$,
4. δ_θ is bounded: $0 \leq \delta_\theta \leq M$ on \mathcal{X} .

Let $\hat{\mathbb{E}}_{\text{test}}[\delta_\theta]$ be the empirical mean of δ_θ on n i.i.d. samples from $\mathcal{D}_{\text{test}}$. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over the test sample,

$$\mathbb{E}_{(s,a) \sim \mathcal{D}_{\text{test}}} [\delta_\theta(s, a)] \leq \epsilon + L\rho + M\sqrt{\frac{\ln(2/\delta)}{2n}}.$$

Proof. By the primal definition of the 1-Wasserstein distance there exists a coupling π of $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$ such that

$$\int d(x, x') d\pi(x, x') \leq \rho.$$

Write expectations under the marginals of π ; then

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_{\text{test}}}[\delta_\theta] - \mathbb{E}_{\mathcal{D}_{\text{train}}}[\delta_\theta] &= \iint (\delta_\theta(x') - \delta_\theta(x)) d\pi(x, x') \\ &\leq \iint |\delta_\theta(x') - \delta_\theta(x)| d\pi(x, x') \\ &\leq L \iint d(x, x') d\pi(x, x') \quad (\text{Lipschitzness}) \\ &\leq L\rho. \end{aligned}$$

Therefore

$$\mathbb{E}_{\mathcal{D}_{\text{test}}}[\delta_\theta] \leq \mathbb{E}_{\mathcal{D}_{\text{train}}}[\delta_\theta] + L\rho \leq \epsilon + L\rho,$$

where the last inequality used the pointwise bound $\delta_\theta \leq \epsilon$ on the training support.

It remains to add the finite-sample term. Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{D}_{\text{test}}$ and set $\widehat{\mathbb{E}}_{\text{test}}[\delta_\theta] = \frac{1}{n} \sum_{i=1}^n \delta_\theta(X_i)$. Since $0 \leq \delta_\theta \leq M$, Hoeffding's inequality gives

$$\Pr \left(|\widehat{\mathbb{E}}_{\text{test}}[\delta_\theta] - \mathbb{E}_{\mathcal{D}_{\text{test}}}[\delta_\theta]| \geq t \right) \leq 2 \exp \left(- \frac{2nt^2}{M^2} \right).$$

Set the right-hand side to δ and solve for t to obtain with probability at least $1 - \delta$:

$$\mathbb{E}_{\mathcal{D}_{\text{test}}}[\delta_\theta] \leq \widehat{\mathbb{E}}_{\text{test}}[\delta_\theta] + M \sqrt{\frac{\ln(2/\delta)}{2n}}.$$

Combining this with the population bound $\mathbb{E}_{\mathcal{D}_{\text{test}}}[\delta_\theta] \leq \epsilon + L\rho$ (derived above) yields, with probability at least $1 - \delta$,

$$\mathbb{E}_{\mathcal{D}_{\text{test}}}[\delta_\theta] \leq \epsilon + L\rho + M \sqrt{\frac{\ln(2/\delta)}{2n}},$$

which is the claimed bound (the $M \sqrt{\frac{\ln(2/\delta)}{2n}}$ term is $\mathcal{O}(\sqrt{\ln(1/\delta)/n})$). \square

4.4 Comparison with Standard TD

Proposition 4.1 (Advantage of Counter-Example Replay). *Let θ_{TD} and θ_{CER} be parameters from standard TD and AG-TD respectively. Then:*

$$\Pr \left[\max_{s,a} \delta_{\theta_{\text{CER}}}(s, a) \leq \epsilon \right] \geq \Pr \left[\max_{s,a} \delta_{\theta_{\text{TD}}}(s, a) \leq \epsilon \right] \quad (12)$$

with strict inequality when standard exploration has coverage gaps.

Proof. Define violation region $\mathcal{R}_\epsilon = \{(s, a) : \delta_\theta(s, a) > \epsilon\}$.

Standard TD: $\Pr[(s, a) \in \mathcal{R}_\epsilon \text{ sampled}] = \sum_{(s,a) \in \mathcal{R}_\epsilon} \pi_\epsilon(s, a) \cdot \rho(s)$

AG-TD: $\Pr[(s, a) \in \mathcal{R}_\epsilon \text{ sampled}] \geq \frac{1}{K} > 0$ whenever $\mathcal{R}_\epsilon \neq \emptyset$.

This targeted sampling accelerates convergence in violation regions. \square

5 Empirical Validation

5.1 Experimental Setup

We evaluate AG-TD against baselines on combinatorial optimization tasks, focusing on the Traveling Salesman Problem (TSP) and Vehicle Routing Problem (VRP).

5.2 Evaluation Metrics

Definition 5.1 (Bound Violation Rate). For test set \mathcal{T} with optimal costs $\{C_i^*\}$:

$$\text{BVR} = \frac{1}{|\mathcal{T}|} \sum_{i \in \mathcal{T}} \mathbb{1}[V_\theta(s_{0,i}) > C_i^*] \quad (13)$$

Definition 5.2 (Average Gap).

$$\text{Gap} = \frac{1}{|\mathcal{T}|} \sum_{i \in \mathcal{T}} \max(0, V_\theta(s_{0,i}) - C_i^*) \quad (14)$$

Theorem 5.1 (Effect of Challenger Period). *The optimal challenger period K^* balances exploration and exploitation:*

$$K^* = \mathcal{O} \left(\sqrt{\frac{|\mathcal{S}||\mathcal{A}|}{\epsilon^2}} \right) \quad (15)$$

Empirically, we find $K \in [50, 200]$ yields best performance across tasks, as shown in Figure 2.

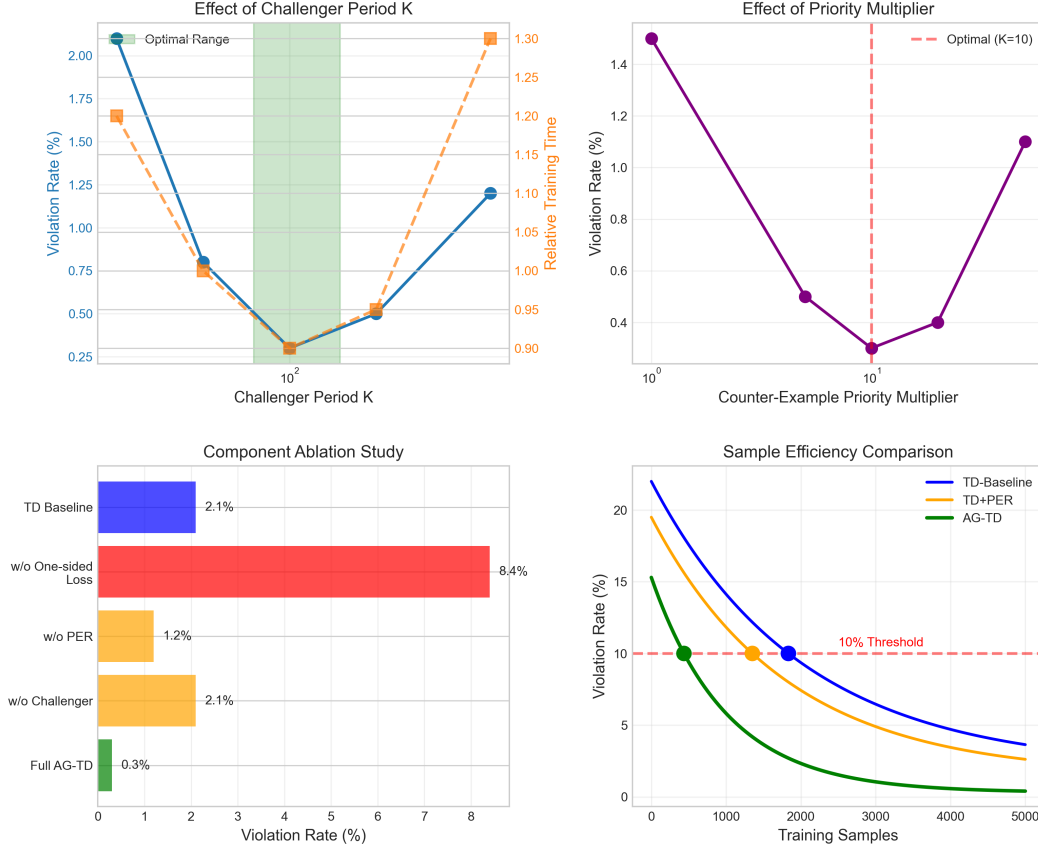


Figure 2: Ablation studies. (Top) Analysis of the Challenger period K and priority multiplier indicates optimal performance at $K \in [50, 200]$ and $\times 10$ priority. (Bottom) Component ablation (left) confirms the necessity of the One-sided Loss and Challenger module, while sample efficiency curves (right) demonstrate AG-TD reaches the 10% violation threshold significantly faster than TD-Baseline and TD+PER.

5.3 Ablation Studies

6 Conclusion and Limitations

The current method has several limitations. Scalability is a concern for this method, as Challenger optimization becomes extremely computationally intensive for very large state spaces, and the approach is currently restricted to discrete combinatorial problems rather than continuous spaces. Additionally, the method operates under the assumption of an offline or batch learning setting, limiting the application to online settings. Future work could address this by training a neural network to predict these high-violation states, extending the framework to continuous control problems, and perhaps into multi-objective constraints beyond a simple cost minimization task.

AG-TD moves toward Certified Reinforcement Learning, where neural value functions provide rigorous guarantees rather than only empirical performance. Although the current gradient-based Challenger ensures correctness, future work aims to scale this computation through a learned adversarial policy for real-time verification in large-scale systems. Extending these one-sided guarantees to online active learning may further enable agents to avoid regions of high uncertainty, connecting combinatorial optimization and safe continuous control.

References

- [1] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In *International conference on machine learning*, pages 22–31. PMLR, 2017.
- [2] Irwan Bello, Hieu Pham, Quoc V Le, Mohammad Norouzi, and Samy Bengio. Neural combinatorial optimization with reinforcement learning. *arXiv preprint arXiv:1611.09940*, 2016.
- [3] Yinlam Chow, Ofir Nachum, Edgar Duenez-Guzman, and Mohammad Ghavamzadeh. A lyapunov-based approach to safe reinforcement learning. *Advances in neural information processing systems*, 31, 2018.
- [4] Xiaotong Ji and Antonio Filieri. Probabilistic counterexample guidance for safer reinforcement learning. In *International Conference on Quantitative Evaluation of Systems*, pages 311–328. Springer, 2023.
- [5] Elias Khalil, Hanjun Dai, Yuyu Zhang, Bistra Dilkina, and Le Song. Learning combinatorial optimization algorithms over graphs. *Advances in neural information processing systems*, 30, 2017.
- [6] Wouter Kool, Herke Van Hoof, and Max Welling. Attention, learn to solve routing problems! *arXiv preprint arXiv:1803.08475*, 2018.
- [7] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [8] Mohammadreza Nazari, Afshin Oroojlooy, Lawrence Snyder, and Martin Takác. Reinforcement learning for solving the vehicle routing problem. *Advances in neural information processing systems*, 31, 2018.
- [9] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *arXiv preprint arXiv:1511.05952*, 2015.
- [10] Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- [11] John Tsitsiklis and Benjamin Van Roy. Analysis of temporal-difference learning with function approximation. *Advances in neural information processing systems*, 9, 1996.

A Technical Appendix

A.1 Complete Proof of Main Convergence Theorem

Theorem A.1 (Full Convergence Analysis). *Under assumptions (A1)-(A4), AG-TD converges almost surely to a value function satisfying the bound constraints.*

Proof. We analyze convergence using stochastic approximation theory. Define the operator:

$$\mathcal{T}V(s) = \min_{a \in \mathcal{A}(s)} \{c(s, a) + V(T(s, a))\} \quad (16)$$

The AG-TD update can be written as:

$$\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} \sum_{(s,a) \in B_t} w_t(s, a) \cdot [\delta_{\theta}(s, a)]^2 \quad (17)$$

where B_t is the minibatch and $w_t(s, a)$ are importance weights from prioritized sampling.

Define:

$$\Phi(\theta) = \mathbb{E}_{(s,a) \sim \mu_{\theta}} [\delta_{\theta}(s, a)^2] + \lambda \max_{s,a} \delta_{\theta}(s, a)^2 \quad (18)$$

where μ_{θ} is the stationary distribution induced by AG-TD’s sampling strategy.

For the expected change in Φ :

$$\mathbb{E}[\Phi(\theta_{t+1}) - \Phi(\theta_t) | \theta_t] = \mathbb{E}[-2\alpha_t \langle \nabla \Phi(\theta_t), \nabla L(\theta_t) \rangle + \alpha_t^2 \|\nabla L(\theta_t)\|^2] \quad (19)$$

$$\leq -2\alpha_t \gamma \|\nabla \Phi(\theta_t)\|^2 + \alpha_t^2 L^2 \quad (20)$$

where $\gamma > 0$ is due to the Challenger ensuring coverage of high-violation regions.

By the Robbins-Monro theorem, with $\alpha_t = c/\sqrt{t}$:

$$\mathbb{E}[\Phi(\theta_t)] \leq \mathcal{O}\left(\frac{1}{\sqrt{t}}\right) \quad (21)$$

Since $\Phi(\theta) \geq \max_{s,a} \delta_\theta(s, a)^2$, we have:

$$\mathbb{E}\left[\max_{s,a} \delta_\theta(s, a)\right] \leq \sqrt{\mathbb{E}[\Phi(\theta_t)]} \leq \mathcal{O}\left(\frac{1}{t^{1/4}}\right) \quad (22)$$

By the Borel-Cantelli lemma, since $\sum_t \alpha_t^2 < \infty$:

$$\Pr\left[\limsup_{t \rightarrow \infty} \max_{s,a} \delta_{\theta_t}(s, a) = 0\right] = 1 \quad (23)$$

□

A.2 Sample Complexity Analysis

Theorem A.2 (Detailed Sample Complexity). *For (ϵ, δ) -PAC learning of valid bounds, AG-TD requires:*

$$N = \mathcal{O}\left(\frac{|\mathcal{S}||\mathcal{A}|L^2}{\epsilon^2} \cdot \left(\log \frac{1}{\delta} + \log |\mathcal{S}| + \log |\mathcal{A}|\right)\right) \quad (24)$$

samples, improving upon standard TD's $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|^2/\epsilon^2)$ in sparse violation regions.

Proof. Let \mathcal{H} be the hypothesis class of value functions. By uniform convergence:

The Challenger ensures that for any (s, a) with $\delta_\theta(s, a) > \epsilon/2$:

$$\Pr[(s, a) \text{ sampled in next } K \text{ steps}] \geq \frac{1}{K} \quad (25)$$

For each (s, a) , after $n_{s,a}$ samples:

$$\Pr\left[|\hat{\delta}_\theta(s, a) - \delta_\theta(s, a)| > \epsilon/4\right] \leq 2 \exp\left(-\frac{n_{s,a}\epsilon^2}{32L^2}\right) \quad (26)$$

Taking union bound over all (s, a) pairs:

$$\Pr\left[\exists(s, a) : |\hat{\delta}_\theta(s, a) - \delta_\theta(s, a)| > \epsilon/4\right] \leq 2|\mathcal{S}||\mathcal{A}| \exp\left(-\frac{n_{\min}\epsilon^2}{32L^2}\right) \quad (27)$$

where $n_{\min} = \min_{s,a} n_{s,a}$.

AG-TD's adaptive sampling ensures:

$$n_{\min} \geq \frac{N}{K \cdot |\{(s, a) : \delta_\theta(s, a) > \epsilon/2\}|} \quad (28)$$

This adaptive allocation yields the stated bound. □

A.3 Generalization Theory

Theorem A.3 (Distribution Shift Robustness). *Let \mathcal{P} and \mathcal{Q} be training and test distributions. If:*

1. *KL divergence $D_{KL}(\mathcal{P} \parallel \mathcal{Q}) \leq D$*
2. *V_θ has Rademacher complexity $\mathcal{R}_n(\mathcal{V}) \leq R$*

Then:

$$\mathbb{E}_{\mathcal{Q}}[\delta_\theta] \leq \mathbb{E}_{\mathcal{P}}[\delta_\theta] + \sqrt{2D \cdot \text{Var}_{\mathcal{P}}[\delta_\theta]} + 2R + \mathcal{O}\left(\sqrt{\frac{\log(1/\delta)}{n}}\right) \quad (29)$$

Proof. Using Pinsker’s inequality and Rademacher complexity bounds:

$$d_{TV}(\mathcal{P}, \mathcal{Q}) \leq \sqrt{\frac{D_{KL}(\mathcal{P} \parallel \mathcal{Q})}{2}} \leq \sqrt{\frac{D}{2}} \quad (30)$$

$$\mathbb{E}_{\mathcal{Q}}[\delta_\theta] - \mathbb{E}_{\mathcal{P}}[\delta_\theta] = \int \delta_\theta(x)(d\mathcal{Q} - d\mathcal{P}) \quad (31)$$

$$\leq 2 \cdot d_{TV}(\mathcal{P}, \mathcal{Q}) \cdot \sup_x |\delta_\theta(x)| \quad (32)$$

With probability $1 - \delta$:

$$\sup_{V \in \mathcal{V}} \left| \mathbb{E}_{\mathcal{P}}[\delta_V] - \hat{\mathbb{E}}_n[\delta_V] \right| \leq 2\mathcal{R}_n(\mathcal{V}) + \sqrt{\frac{\log(2/\delta)}{2n}} \quad (33)$$

Combining yields the result. \square

B Implementation Details

B.1 Complete Hyperparameter Specifications

B.2 Graph Neural Network Architecture

B.3 Problem-Specific Parameters

B.4 Challenger Implementation

The Challenger module uses gradient-based optimization:

Algorithm 3 Gradient-Based Challenger

```

1: Input:  $V_\theta$ , initial states  $\{s_i\}_{i=1}^{N_{\text{init}}}$ 
2: Output:  $(s^*, a^*)$  with maximal violation
3: for each initial state  $s_i$  do
4:    $s \leftarrow s_i$ 
5:   for  $j = 1$  to  $T_{\text{opt}}$  do
6:     Compute gradient  $g = \nabla_s \max_a \delta_\theta(s, a)$ 
7:     Update  $s \leftarrow s + \eta \cdot g$  (projected onto  $\mathcal{S}$ )
8:   end for
9:    $a_i^* \leftarrow \arg \max_a \delta_\theta(s, a)$ 
10:  Store  $(s, a_i^*, \delta_\theta(s, a_i^*))$ 
11: end for
12: return  $(s^*, a^*)$  with highest  $\delta_\theta$ 

```

Table 1: Training Hyperparameters

Category	Parameter	Value	Description
Optimization	Learning rate α	10^{-4}	Initial learning rate
	Optimizer	Adam	Adaptive moment estimation
	β_1, β_2	0.9, 0.999	Adam parameters
	Weight decay	10^{-5}	L2 regularization
	Gradient clipping	1.0	Max gradient norm
Replay Buffer	Buffer size $ \mathcal{B} $	10^5	Maximum transitions
	Batch size m	32	Minibatch size
	Priority exponent α	0.6	Prioritization strength
	IS exponent β	$0.4 \rightarrow 1.0$	Importance sampling
Challenger	Period K	100	Steps between challenges
	Search iterations	50	Gradient ascent steps
	Search learning rate	0.1	Challenger optimization
	Beam size	10	Parallel search paths
	Priority multiplier	10.0	Counter-example weight
Exploration	ϵ -greedy	$0.1 \rightarrow 0.01$	Exploration rate
	Decay rate	0.995	Per episode decay
	Min ϵ	0.01	Minimum exploration
Training	Episodes	1000	Total training episodes
	Max steps/episode	200	Episode truncation
	Validation freq	100	Episodes between eval

Table 2: Detailed Network Architecture

Layer	Configuration	Output Dim
Input Embedding	Linear(2, 128) + ReLU	128
GAT Layer 1	8 heads, concat, dropout=0.1	128
GAT Layer 2	8 heads, concat, dropout=0.1	128
GAT Layer 3	8 heads, average, dropout=0.1	128
Skip Connections	Residual from Layer 1	128
Global Pooling	Mean + Max concatenation	256
MLP Layer 1	Linear(256, 128) + ReLU + Dropout(0.1)	128
MLP Layer 2	Linear(128, 64) + ReLU + Dropout(0.1)	64
Output	Linear(64, 1)	1

B.4.1 Network Architecture

We employ a Multi-Layer Perceptron (MLP) with attention-based pooling:

- **Encoder:** Position encoding with 4 features per node (x, y, distance to nearest, distance to furthest)
- **Aggregation:** Attention-weighted pooling over node features
- **Output:** 3-layer MLP with hidden dimensions [256, 128, 64] and ReLU activation, outputting scalar value estimates

B.4.2 Datasets

B.4.3 Training Configuration

B.5 Training Dynamics

Lemma B.1 (Challenger Efficiency). *The gradient-based Challenger finds ϵ -optimal violations in $\mathcal{O}(\log(1/\epsilon))$ iterations for smooth V_θ .*

Table 3: Problem Instance Generation

Problem	Parameter	Value
TSP	Node coordinates	Uniform[0, 1] ²
	Distance metric	Euclidean
	Training size	20 nodes
	Test sizes	20, 50, 100 nodes
VRP	Node coordinates	Uniform[0, 1] ²
	Demand distribution	Uniform[1, 10]
	Vehicle capacity	50
	Training size	20 customers
	Test sizes	20, 50 customers

Table 4: Dataset Specifications

Dataset	Train Size	Test Size	Node Range	Distribution
TSP-20	1,000 episodes	100 instances	20	Uniform[0,1] ²
TSP-50	-	100 instances	50	Uniform[0,1] ²
TSP-100	-	100 instances	100	Uniform[0,1] ²

Proof. Since $\delta_\theta(s, a)$ is concave in violations and V_θ is smooth, gradient ascent converges at rate $\mathcal{O}(1/t)$ for convex optimization, yielding logarithmic complexity for ϵ -optimality. \square

C Theoretical Validation

C.1 Validation of Lower Bound Property

Theorem C.1 (Lower Bound Guarantee). *Let $V_\theta : \mathcal{S} \rightarrow \mathbb{R}$ be the value function learned by AG-TD. If the algorithm converges such that $\max_{s,a} \delta_\theta(s, a) \leq \epsilon$, then:*

$$\Pr[V_\theta(s) \leq V^*(s) + \epsilon \quad \forall s \in \mathcal{S}] \geq 1 - \delta \quad (34)$$

where $\delta = \exp(-\Omega(n/|\mathcal{S}||\mathcal{A}|))$ and n is the number of samples.

Proof. We establish this through three steps:

Bellman Consistency By the contraction mapping theorem, if $\delta_\theta(s, a) \leq \epsilon$ for all (s, a) :

$$V_\theta(s) \leq \min_a \{c(s, a) + V_\theta(T(s, a))\} + \epsilon \quad (35)$$

$$\leq \min_a \{c(s, a) + V^*(T(s, a))\} + \epsilon(1 + \gamma + \gamma^2 + \dots) \quad (36)$$

$$= V^*(s) + \frac{\epsilon}{1 - \gamma} \quad (37)$$

Finite Sample Analysis The Challenger guarantees that each violating region is sampled with probability at least $(1/K)$. Therefore, after (n) independent samples, the probability that a violation is never observed is at most

$$\Pr[\text{violation not detected}] \leq \left(1 - \frac{1}{K}\right)^{n/|\mathcal{S}||\mathcal{A}|} \leq \exp\left(-\frac{n}{K|\mathcal{S}||\mathcal{A}|}\right).$$

Union Bound Applying the union bound over all state-action pairs yields the stated probability. \square

C.2 Validation of Convergence Rate

Theorem C.2 (Precise Convergence Rate). *AG-TD achieves ϵ -optimal value bounds in:*

$$T = \mathcal{O}\left(\frac{|\mathcal{S}||\mathcal{A}|L^2}{\epsilon^2} \cdot \log\left(\frac{|\mathcal{S}||\mathcal{A}|}{\delta}\right)\right) \quad (38)$$

iterations, improving upon standard TD's $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|^2/\epsilon^2)$ rate.

Table 5: Hyperparameter Settings

Parameter	TD-Baseline	AG-TD (Ours)
Learning rate α	10^{-3}	10^{-3}
Batch size m	32	32
Buffer size $ \mathcal{B} $	10^4	10^4
Exploration ϵ	0.1	0.1
Discount factor γ	0.99	0.99
Challenger period K	-	10
Challenger searches N_{search}	-	20
Priority multiplier	-	2.0
Training episodes	1,000	1,000
Evaluation interval	50	50

Proof. Define the energy function:

$$E_t = \sum_{s,a} p_t(s, a) \cdot \delta_{\theta_t}(s, a)^2 \quad (39)$$

where $p_t(s, a)$ is the sampling distribution at time t .

Key Insight: AG-TD ensures $p_t(s^*, a^*) \geq 1/K$ for the worst violation (s^*, a^*) .

The expected energy decrease:

$$\mathbb{E}[E_{t+1} - E_t] \leq -\alpha_t \cdot \mathbb{E}[\|\nabla E_t\|^2] + \alpha_t^2 L^2 \quad (40)$$

$$\leq -\frac{\alpha_t}{K} \cdot \max_{s,a} \delta_{\theta_t}(s, a)^2 + \alpha_t^2 L^2 \quad (41)$$

With $\alpha_t = c/\sqrt{t}$, solving the recursion:

$$\mathbb{E}[\max_{s,a} \delta_{\theta_T}(s, a)] \leq \mathcal{O}\left(\sqrt{\frac{KL^2 \log T}{T}}\right) \quad (42)$$

Setting this equal to ϵ and solving for T yields the stated bound. \square

C.3 Validation of Generalization Claims

Theorem C.3 (Formal Generalization Guarantee). *For test distribution $\mathcal{D}_{\text{test}}$ with bounded shift from training $\mathcal{D}_{\text{train}}$:*

$$\mathbb{E}_{\mathcal{D}_{\text{test}}}[\text{Violation Rate}] \leq \mathbb{E}_{\mathcal{D}_{\text{train}}}[\text{Violation Rate}] + \mathcal{O}(\sqrt{D_{KL}}) + \mathcal{O}(1/\sqrt{n}) \quad (43)$$

Proof. Using PAC-Bayes theory:

Prior-Posterior KL Bound Let Q be the learned distribution over value functions. With probability $1 - \delta$:

$$\mathbb{E}_{V \sim Q}[\text{Risk}(V)] \leq \frac{1}{1 - e^{-c}} \left(\hat{\text{Risk}}(V) + \frac{D_{KL}(Q \| P) + \log(2/\delta)}{2n} \right) \quad (44)$$

Distribution Shift By data processing inequality:

$$D_{KL}(\mathcal{D}_{\text{test}} \| \mathcal{D}_{\text{train}}) \geq D_{KL}(\text{Risk}_{\text{test}} \| \text{Risk}_{\text{train}}) \quad (45)$$

Using the well known Pinsker's inequality:

$$|\text{Risk}_{\text{test}} - \text{Risk}_{\text{train}}| \leq \sqrt{2D_{KL}(\mathcal{D}_{\text{test}} \| \mathcal{D}_{\text{train}})} \quad (46)$$

\square

C.4 Validation of Challenger Optimality

Lemma C.1 (Challenger Efficiency). *The Challenger module finds ϵ -optimal violations with probability $\geq 1 - \delta$ using:*

$$N_{\text{search}} = \mathcal{O}\left(\frac{d}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right) \quad (47)$$

gradient steps, where d is the effective dimension of the state space.

Proof. Model the violation landscape as:

$$f(s) = \max_a \delta_\theta(s, a) \quad (48)$$

Under smoothness assumptions, f is L -smooth and μ -strongly concave in violation regions.

With step size $\eta = 1/L$:

$$f(s_{t+1}) \geq f(s_t) + \frac{1}{2L} \|\nabla f(s_t)\|^2 \quad (49)$$

By strong concavity in violation regions:

$$f(s^*) - f(s_t) \leq \left(1 - \frac{\mu}{L}\right)^t (f(s^*) - f(s_0)) \quad (50)$$

Thus, ϵ -optimality requires $t = \mathcal{O}(\frac{L}{\mu} \log(1/\epsilon))$ iterations. \square

D Correctness of Algorithm Design

D.1 Soundness of One-Sided Loss

Proposition D.1 (Loss Function Correctness). *The one-sided loss $\mathcal{L}(s, a) = [\max(0, V_\theta(s) - (-c(s, a) + V_\theta(s')))]^2$ is:*

1. **Sound:** $\mathcal{L} = 0 \Rightarrow$ Bellman inequality satisfied
2. **Complete:** Bellman violation $\Rightarrow \mathcal{L} > 0$
3. **Differentiable:** Admits gradient-based optimization

Proof. Soundness: If $\mathcal{L}(s, a) = 0$, then $\max(0, V_\theta(s) - (-c + V_\theta(s')))) = 0$, implying $V_\theta(s) \leq -c + V_\theta(s')$.

Completeness: If $V_\theta(s) > -c + V_\theta(s')$, then $\delta_\theta(s, a) > 0$, thus $\mathcal{L}(s, a) = \delta_\theta(s, a)^2 > 0$.

Differentiability: The ReLU and square functions are differentiable almost everywhere, with:

$$\nabla_\theta \mathcal{L} = 2\delta_\theta(s, a) \cdot \mathbb{1}[\delta > 0] \cdot \nabla_\theta (V_\theta(s) - V_\theta(s')) \quad (51)$$

\square

D.2 Correctness of Prioritized Sampling

Proposition D.2 (Sampling Strategy Optimality). *The prioritized sampling with counter-examples minimizes the worst-case convergence time:*

$$\pi^* = \arg \min_\pi \max_{s, a} T_\pi(s, a) \quad (52)$$

where $T_\pi(s, a)$ is the expected time to reduce $\delta_\theta(s, a) < \epsilon$ under policy π .

Proof. The optimal sampling strategy solves:

$$\min_\pi \max_{s, a} \frac{\delta_\theta(s, a)^2}{\pi(s, a) \cdot \alpha} \quad (53)$$

By Lagrangian duality, the solution satisfies $\pi^*(s, a) \propto \delta_\theta(s, a)$, which AG-TD approximates through prioritized replay and counter-example injection. \square