

---

# A Near-Optimal Control Policy for Data-driven Assemble-to-order Systems

---

**Lun Yu**  
CUHK-Shenzhen  
yulun@cuhk.edu.cn

**Zhixuan Cai**  
Tsinghua University  
caizx21@mails.tsinghua.edu.cn

**Zhaoran Wang**  
Northwestern University  
zhaoranwang@gmail.com

**Tianhu Deng**  
Soochow University  
thdeng@suda.edu.cn

## Abstract

We study a data-driven assemble-to-order (ATO) control problem, aiming to synchronize component ordering and product assembly under unknown demand distributions and non-identical lead times. We address two key questions: the statistical tractability of learning a near-optimal policy from limited data and the computational complexity of obtaining it. To the best of our knowledge, our work is the first to analyze the sample efficiency for a general ATO system as a multidimensional control problem and to propose an algorithm that finds a provably near-optimal solution. Methodologically, we introduce a novel asymmetric Lipschitz continuity (ALC) property to establish regularity conditions for the infinite-horizon problem. Surprisingly, we prove that the data-driven ATO problem avoids the curse of dimensionality; the performance gap of our policy scales as  $O(M^{-1/2} \log M)$  with sample size  $M$ , only logarithmically worse than approximating the demand mean. We develop a specialized reinforcement learning (RL) algorithm that exploits a convex-preserving property in ATO dynamics, using input convex neural networks and interior point methods to achieve computational feasibility. Numerical studies show our algorithm consistently and significantly outperforms existing heuristics and a general-purpose RL benchmark.

## 1 Introduction

Assemble-to-order (ATO) is a widespread production strategy where products are assembled from components only after customer orders are received. This allows companies like Dell and BMW [9, 12] to offer high product variety while minimizing finished-goods inventory. However, managing ATO systems is challenging due to the need to coordinate procurement of numerous components with varying lead times to meet uncertain future demand for multiple products, and is particularly challenging when the demand distribution is unknown and must be learned from historical data.

Motivated by our collaboration with a processed-food manufacturer, which delivers a variety of 8 products using 15 components, we consider the ATO control problem in a data-driven setting. We formulate it as an infinite-horizon Markov Decision Process (MDP) and address two fundamental questions:

1. **Statistical Tractability:** How does the performance of the best policy learned from a finite demand sample of size  $M$  scale with the problem dimension? Can we learn a near-optimal policy without an impractically large dataset?

2. **Computational Complexity:** Can we efficiently compute such a policy, given that the optimal policy structure is unknown and the state-action space is high-dimensional?

**Our Contributions.** Our work provides affirmative answers to these questions.

First, we show that a Convex Fitted Q-Iteration (CFQI) algorithm efficiently computes a near-optimal policy. The optimality gap of CFQI is  $O(p^2 M^{-1/2} \log M)$  relative to the intractable optimal policy, where  $p$  is the dimension of the system. This demonstrates that accurately approximating the high-dimensional optimal value function requires only logarithmically more samples than approximating the demand mean, which avoids the curse of dimensionality. Our theoretical analysis introduces a novel *asymmetric Lipschitz continuity (ALC)* property. Standard Lipschitz conditions are often too restrictive for infinite-horizon inventory problems. ALC provides a weaker, more suitable regularity condition that reflects the inherent asymmetry between over-stocking and under-stocking costs, enabling a rigorous sample complexity analysis for the ATO control problem.

Second, we propose a specialized reinforcement learning (RL) algorithm as implementation of CFQI. We leverage a convex-preserving property of the control problem by employing Input Convex Neural Networks (ICNNs) [2] to approximate the convex value functions, transforming the complex policy improvement step into a series of unconstrained convex optimizations. Our extensive numerical results validate our theory, showing that our algorithm significantly outperforms strong baselines, including a general-purpose TD3 algorithm and established heuristics.

## 2 Data-driven ATO as a Markov Decision Process

We model an ATO system with  $m$  products and  $n$  components. The assembly structure is defined by a Bill-of-Materials (BOM) matrix  $A \in \mathbb{R}^{m \times n}$ . Component  $j$  has a lead time of  $L_j \geq 1$  periods. In each period  $t$ , the system state  $\mathbf{s}(t) = (\mathbf{x}(t), \tilde{\mathbf{u}}(t), R(t))$  is observed, where  $\mathbf{x}(t) \in \mathbb{R}_+^n$  is the on-hand component inventory,  $\tilde{\mathbf{u}}(t) \in \mathbb{R}_+^m$  is the total backlog (including new demand  $\mathbf{d}(t)$ ), and  $R(t) \in \mathbb{R}^{n \times L_n}$  is the pipeline inventory matrix. The demand vector  $\mathbf{d}(t)$  is i.i.d. from an unknown distribution  $F$  with compact support, but we have access to a historical sample  $\hat{D} = \{\mathbf{d}_1, \dots, \mathbf{d}_M\}$ .

The actions are  $\mathbf{a}(t) = (\mathbf{y}(t), \mathbf{z}(t+1))$ , where  $\mathbf{y}(t) \in \mathbb{R}_+^m$  is the assembly quantity and  $\mathbf{z}(t+1) \in \mathbb{R}_+^n$  is the component order quantity for the next cycle. Actions are constrained by available inventory:  $A^\top \mathbf{y}(t) \leq \mathbf{x}(t)$  and  $\mathbf{y}(t) \leq \tilde{\mathbf{u}}(t)$ . The single-period cost is  $C(\mathbf{s}(t), \mathbf{a}(t)) = \mathbf{h}^\top (\mathbf{x}(t) - A^\top \mathbf{y}(t)) + \mathbf{b}^\top (\tilde{\mathbf{u}}(t) - \mathbf{y}(t))$ , where  $\mathbf{h}$  and  $\mathbf{b}$  are component holding and product backorder costs, respectively. The system evolves according to the transition function  $\mathbf{s}(t+1) = f(\mathbf{s}(t), \mathbf{a}(t), \mathbf{d}(t+1))$ . The objective is to find a policy  $\pi$  that minimizes the total expected discounted cost  $V^\pi(\mathbf{s}; F) = \mathbb{E} [\sum_{t=0}^{\infty} \gamma^t C(\mathbf{s}(t), \mathbf{a}(t))]$ , where the expectation is over future demands from the true distribution  $F$ . The optimal value function is  $V^*(\mathbf{s}; F) = \inf_{\pi} V^\pi(\mathbf{s}; F)$ . For more details on the ATO model, please refer to Appendix A.

## 3 Algorithm and Theoretical Guarantees

Directly solving the MDP is intractable due to the unknown distribution  $F$  and the continuous, high-dimensional state-action spaces. We propose a solution based on *Sample Average Approximation (SAA)* combined with a specialized RL algorithm. To handle the unbounded state space, we first introduce a *Damped Sample Average Approximation (DSAA)* of the problem. We introduce a damping factor  $\alpha \in (0, 1)$  to the state transition, making the state and action spaces provably bounded while preserving the linearity of the dynamics and convexity of the state-action space. We then replace the unknown demand distribution  $F$  with the empirical distribution  $\hat{F}$  from the sample  $\hat{D}$ . This defines a tractable MDP approximation with a known transition kernel.

The optimal action-value function  $Q^*$  for this DSAA problem can be found as the unique fixed point of the Bellman optimality operator,  $\mathcal{T}^*Q = Q$ , where for a state-action pair  $(\mathbf{s}, \mathbf{a})$ :

$$(\mathcal{T}^*w)(\mathbf{s}, \mathbf{a}) := C(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\hat{\mathbf{d}} \sim \hat{F}} \left[ \min_{\mathbf{a}' \in \mathcal{A}_{\alpha}(\mathbf{s}')} w(\mathbf{s}', \mathbf{a}') \right], \quad \text{for } \mathbf{s}' = f_{\alpha}(\mathbf{s}, \mathbf{a}, \hat{\mathbf{d}}). \quad (1)$$

Our first key insight is that this operator preserves convexity.

**Theorem 1 (Convexity Preservation).** *If the function  $w$  is convex over the state-action space, then  $\mathcal{T}^*w$  is also convex.*

This theorem implies that the optimal action-value function  $Q^*$  of the DSAA problem is convex. This property is critical, as it allows us to restrict the search for  $Q^*$  to the space of convex functions, which can be efficiently approximated and optimized.

Based on the convexity property, we propose the Convex Fitted Q-Iteration (CFQI) algorithm (Algorithm 1). CFQI approximates the action-value function  $Q^*$  by iteratively applying the Bellman operator to a sequence of convex functions  $\{Q_k\}$ . In practice, we parameterize  $Q_k$  using *Input Convex Neural Networks (ICNNs)*, a neural network architecture whose output is guaranteed to be a convex function of its input [2].

The policy improvement step,  $\min_{\mathbf{a}} Q_k(\mathbf{s}, \mathbf{a})$ , is a constrained convex optimization problem. To handle this efficiently across thousands of states in each iteration, we employ a penalty method, transforming it into an unconstrained convex optimization that can be solved rapidly using gradient-based methods. The final output is a policy  $\hat{\pi}^*(\hat{D})$  derived from the learned convex Q-function, translated back to be applicable to the original, non-damped system.

---

**Algorithm 1:** Convex FQI.

---

**Input:** Number of iterations  $K$ , Class  $\mathcal{G}$  with accuracy level  $\varepsilon$ , regressor **Regress**, initial actor-value function  $Q_0 \in \mathcal{C} \cap \mathcal{G}$ , sample  $\hat{D}$ , initial state  $\mathbf{s} = \mathbf{s}(0)$ , and constants  $\alpha \in (0, \bar{d}/\|\mathbf{s}\|_\infty)$  and  $\eta \in (0, 1)$ .

**Output:** Policy  $\hat{\pi}^*(\hat{D})$ .

**for**  $k = 0$  **to**  $K - 1$  **do**

    apply **Regress** to obtain convex function

$$Q_{k+1} \in \text{Regress}\left(\mathcal{T}^*(Q_k(\cdot) + \Lambda(\cdot; \eta^k))\right), \quad (2)$$

    where the Bellman operator  $\mathcal{T}^*$  is characterized by taking  $G = \hat{F}$ .

take a map  $\xi_K$  on  $\mathcal{S}_\alpha$  satisfying

$$\xi_K(\mathbf{s}) \in \arg \min_{\mathbf{a} \in \mathcal{A}(\mathbf{s})} Q_K(\mathbf{s}, \mathbf{a}).$$

return policy  $\hat{\pi}^*(\hat{D}) := \{\hat{\mathbf{a}}(t) : t \geq 0\}$ , for recursively defined sequences

$$\hat{\mathbf{a}}(t) := (\hat{\mathbf{y}}(t), \hat{\mathbf{z}}(t+1)) := \xi_K(\hat{\mathbf{s}}(t) - \mathbf{r}(t)), \quad (3)$$

$$\hat{\mathbf{s}}(t+1) = f(\hat{\mathbf{s}}(t), \hat{\mathbf{a}}(t), \mathbf{d}(t+1)), \quad \text{and}$$

$$\mathbf{r}(t+1) = (1 - \alpha)\mathbf{r}(t) + \alpha(\mathbf{x}(t) - A^\top \mathbf{y}(t), \tilde{\mathbf{u}}(t) - \mathbf{y}(t), O). \forall t = 0, 1, \dots, \quad (4)$$

with  $\hat{\mathbf{s}}(0) = \mathbf{s}$  and  $\mathbf{r}(0) = \mathbf{0}$ . Note that  $\{\hat{\mathbf{s}}(t)\}$  is the controlled system state.

---

A major challenge in analyzing the sample complexity is establishing regularity of the value function. Standard Lipschitz continuity conditions [8] do not hold for the ATO problem. To overcome this, we introduce a novel property.

**Theorem 2 (Asymmetric Lipschitz Continuity (ALC)).** *The Bellman operator  $\mathcal{T}^*$  preserves a specific form of asymmetric Lipschitz continuity. This implies that the optimal value function  $V_\alpha^*(\cdot; G)$  is uniformly Lipschitz continuous for any demand distribution  $G$ .*

The ALC property intuitively captures the different cost implications of positive vs. negative changes in inventory levels (i.e., overage versus underage). Detailed definition for the ALC property can be found in Appendix B. Armed with this tool, we establish our main theoretical result on the performance of the policy generated by CFQI.

**Theorem 3 (Algorithmic Error).** *For any  $\delta > 0$ ,  $\mathbf{s} \in \mathcal{S}$ , and distribution  $F \in \mathcal{D}$ , the output data driven policy  $\hat{\pi}^*$  from Algorithm 1 with*

$$\alpha = \min \left\{ \frac{1 - \gamma}{M^{1/2}}, \frac{\bar{d}}{\|\mathbf{s}\|_\infty} \right\}, \quad \eta = \max \left\{ \frac{\gamma}{2}, 2\gamma - 1 \right\}, \quad \text{and}$$

$$K \geq 2 + \frac{2 \log \alpha + \log \varepsilon_m - \log(3\bar{a}\bar{c}\bar{d}L_n(m+n))}{\log(1 - \gamma)}$$

gives a data-driven policy  $\hat{\pi}^*$  satisfying

$$V^{\hat{\pi}^*(\hat{D})}(\mathbf{s}; F) - V^*(\mathbf{s}; F) \leq \frac{1}{\sqrt{M}} V^*(\mathbf{s}; F) + \frac{8\varepsilon_m}{(1-\gamma)^3} + \mathcal{E}(\delta) \quad (5)$$

with probability at least  $1 - \delta$ . Here,

$$\mathcal{E}(\delta) = O\left(\frac{\bar{a}\bar{c}(p_S + p_A)^2}{(1-\gamma)^2} \sqrt{\frac{1}{M} \log \frac{\bar{a}M}{\alpha\delta}}\right)$$

and  $O(\cdot)$  hides a constant that does not depend on the model primitives.

Theorem 3 shows that the optimality gap shrinks at a rate of nearly  $M^{-1/2}$  and grows polynomially in dimension, thus avoiding the curse of dimensionality. This guarantees that a near-optimal policy can be learned from a reasonably sized dataset.

## 4 Numerical Studies

We validate our proposed algorithm, CTD3 (a practical implementation of CFQI using TD3 architecture with ICNN critics), against two strong baselines: 1) **PTD3**, a penalized TD3 algorithm using standard (non-convex) networks, and 2) **NV-PRP**, a state-of-the-art heuristic based on newsvendor decomposition and priority rules [6]. We test on three systems of increasing complexity: a small N-system (2 components, 2 products) where the optimal policy is known, a mid-sized PC-system (6 components, 4 products), and a large-scale real-world system from a food manufacturer (15 components, 8 products).

**Results.** Across all settings, CTD3 demonstrates superior performance.

- **N-System:** With independent demands, all methods perform well. However, when demands are negatively correlated (substitutes), CTD3 provides a cost reduction of over 14% compared to the NV-PRP heuristic, which struggles with demand correlation. CTD3 also converges faster and more stably than PTD3.
- **PC-System:** CTD3 consistently achieves lower costs than both PTD3 and NV-PRP across all sample sizes. The performance improves rapidly as  $M$  increases from 10 to 200, confirming our theory that a large sample is not required.
- **Large-System:** In the real-world industrial case (Table 6), the advantage of exploiting convexity becomes crucial. CTD3 achieves a cost of 32.0, an 11.6% improvement over NV-PRP. In contrast, the general-purpose PTD3 algorithm fails to find a competitive policy, highlighting the importance of incorporating problem structure into the RL algorithm design for high-dimensional control tasks.

More detailed results can be found in Appendix E.

**Insights on Heuristics.** Our learned policy reveals insights into common heuristics. In the PC-system, we found that the CTD3 policy frequently violates the well-known "no-holdback" rule, which dictates that available components should always be used to satisfy current demand. By strategically holding back components for more valuable future products, our policy achieves lower costs. This suggests that the optimal no-holdback rule for simpler systems can be significantly suboptimal in more complex, multi-product settings, potentially incurring double-digit performance losses.

## 5 Conclusion

We developed a specialized RL algorithm to solve the data-driven ATO control problem, providing both strong theoretical guarantees and superior empirical performance. Our analysis, enabled by the novel ALC property, is the first to show that this high-dimensional control problem is statistically tractable, avoiding the curse of dimensionality. Our algorithm, by exploiting the problem's inherent convexity with ICNNs, provides a new and effective benchmark for complex ATO systems. This work demonstrates the power of integrating structural properties from operations research into modern deep reinforcement learning frameworks to solve challenging real-world control problems.

## References

- [1] Praveen Agarwal, Mohamed Jleli, and Bessem Samet. *Fixed point theory in metric spaces: Recent advances and applications*. Springer Singapore, Singapore, 2018.
- [2] Brandon Amos, Lei Xu, and J. Zico Kolter. Input convex neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 146–155. PMLR, 06–11 Aug 2017.
- [3] Dimitri Bertsekas and Steven E Shreve. *Stochastic optimal control: The discrete-time case*, volume 5. Athena Scientific, 1996.
- [4] Dimitri P Bertsekas et al. Dynamic programming and optimal control 3rd edition, volume ii. *Belmont, MA: Athena Scientific*, 1, 2011.
- [5] David Blackwell. Discounted dynamic programming. *The Annals of Mathematical Statistics*, 36(1):226–235, 1965.
- [6] Shuyu Chen, Lijian Lu, Jing-Sheng Jeannette Song, and Hanqin Zhang. Optimizing assemble-to-order systems: Decomposition heuristics and scalable algorithms. *SSRN Electronic Journal*, 2021.
- [7] Levi DeValve, Saša Pekeč, and Yehua Wei. A primal-dual approach to analyzing ATO systems. *Management Science*, 66(11):5389–5407, November 2020.
- [8] Karl Hinderer. Lipschitz continuity of value functions in markovian decision processes. *Mathematical Methods of Operations Research*, 62(1):3–22, 2005.
- [9] Roman Kapuscinski, Rachel Q. Zhang, Paul Carbonneau, Robert Moore, and Bill Reeves. Inventory decisions in Dell’s supply chain. *INFORMS Journal on Applied Analytics*, 34(3):191–205, June 2004.
- [10] Lijian Lu, Jing-Sheng Song, and Hanqin Zhang. Optimal and asymptotically optimal policies for assemble-to-order N- and W-systems. *Naval Research Logistics*, 62(8):617–645, 2015.
- [11] Yingdong Lu and Jing-Sheng Song. Order-based cost optimization in assemble-to-order systems. *Operations Research*, 53(1):151–169, February 2005.
- [12] Joann Muller. BMW’s push for made-to-order cars. *Forbes*, 2010. URL <https://www.forbes.com/forbes/2010/0927/companies-bmw-general-motors-cars-bespoke-auto.html>.
- [13] Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research*, 9:815–857, jun 2008.
- [14] Martin J. Wainwright. *High-dimensional statistics: a non-asymptotic viewpoint*. Cambridge University Press, first edition, February 2019.
- [15] Yaqi Xie, Will Ma, and Linwei Xin. Vc theory for inventory policies. *arXiv preprint arXiv:2404.11509*, 2024.

## Appendix

### A MDP Formulation of the ATO Problem

The ATO control problem can be analyzed using an infinite-horizon Markov Decision Process (MDP) formulation. In the formulation, states are updated whenever new information is realized, at which actions are also taken. Therefore, the state of the ATO system takes the form of  $\mathbf{s}(t) := (\mathbf{x}(t), \tilde{\mathbf{u}}(t), R(t))$ , which updates immediately after *observing new demands*. Here

- $\mathbf{x}(t)$  is the inventory level after receiving component replenishment,
- $\tilde{\mathbf{u}}(t) := \mathbf{u}(t) + \mathbf{d}(t)$  is the augmented unsatisfied demands,
- $R(t) \in \mathbb{R}^{n \times L_n}$  is the pipeline inventory matrix after procurement  $\mathbf{z}(t)$  is executed.

For  $j \in [n]$  and  $\ell \in [L_j]$ , the element on the  $(j, \ell)$ -th entry of  $R(t)$  is the quantity of component  $j$  ordered in period  $t + \ell - L_j$ , which is scheduled to arrive in period  $t + \ell$ . The entry is 0 if  $j > L_j$ .

We emphasize that the update of state  $\mathbf{s}(t)$  occurs at a different time from the advancement of  $\mathbf{x}(t)$  to  $\mathbf{x}(t+1)$ . Consequently, we specify the  $t$ -th action by  $\mathbf{a}(t) := (\mathbf{y}(t), \mathbf{z}(t+1))$  since the product assembly quantity  $\mathbf{y}(t)$  and component orders  $\mathbf{z}(t+1)$  are both decided between the observations of demands  $\mathbf{d}(t)$  and  $\mathbf{d}(t+1)$ . Let  $p_S := m + n + nL_n$  and  $p_A := m + n$  be the dimensional of the state and action vector, respectively, the state transition is characterized by a function  $f : \mathbb{R}^{p_S + p_A + m} \rightarrow \mathbb{R}^{p_S}$ ,

$$\mathbf{s}(t+1) = f(\mathbf{s}(t), \mathbf{a}(t), \mathbf{d}(t+1)), \forall t \in \mathbb{Z}_+. \quad (6)$$

Given  $\mathbf{s} = (\mathbf{x}, \tilde{\mathbf{u}}, R)$  and  $\mathbf{a} = (\mathbf{y}, \mathbf{z})$ , the function  $f$  is characterized by  $f(\mathbf{s}, \mathbf{a}, \mathbf{d}) = (\mathbf{x}', \tilde{\mathbf{u}}', R')$ , where  $\mathbf{x}' = \mathbf{x} - A^\top \mathbf{y} + R\mathbf{e}_1$ ,  $\tilde{\mathbf{u}}' = \tilde{\mathbf{u}} + \mathbf{d} - \mathbf{y}$ , and  $R' \in \mathbb{R}^{n \times L_n}$  satisfies that

$$R'_{j,\ell} = \begin{cases} R_{j,\ell+1}, & \text{if } \ell < L_j, \\ z_j, & \text{if } \ell = L_j, \\ 0, & \text{if } \ell > L_j. \end{cases}$$

The manufacturer starts to control the system after the first demand is realized. The manufacturer observes the state  $\mathbf{s}(0)$  and has control over the actions  $\mathbf{a}(t)$  for  $t \in \mathbb{Z}_+$ . We represent the control policy of the ATO system by a stochastic process  $\pi := \{\mathbf{a}(t) : t \in \mathbb{Z}_+\}$ , which encapsulates all procurement and assembly decisions to be made. Let  $\mathcal{R}$  to be the family of  $n \times L_n$  nonnegative matrices  $R$  such that  $R_{j,\ell} = 0$  for  $\ell > L_j$ . Let  $\mathcal{S} := \mathbb{R}_+^{n+m} \times \mathcal{R} \subseteq \mathbb{R}_+^{p_S}$  be the state space. For state  $\mathbf{s} = (\mathbf{x}, \tilde{\mathbf{u}}, R) \in \mathcal{S}$ , denote by  $\mathcal{A}(\mathbf{s})$  the state-dependent action space, i.e.,  $\mathcal{A}(\mathbf{s}) := \{(\mathbf{y}, \mathbf{z}) \in \mathbb{R}_+^{m+n} : A^\top \mathbf{y} \leq \mathbf{x}, \mathbf{y} \leq \tilde{\mathbf{u}}\}$ . Note that  $\mathcal{S}$  is a noncompact subset of  $\mathbb{R}^{p_S}$ . For each  $\mathbf{s} \in \mathcal{S}$ ,  $\mathcal{A}(\mathbf{s})$  is a noncompact subset  $\mathbb{R}^{p_A}$ . Policy  $\pi$  is admissible if: (i) the decision  $\mathbf{a}(t) \in \mathcal{A}(\mathbf{s}(t))$  for all  $t \in \mathbb{Z}_+$ , where the controlled state process  $\{\mathbf{s}(t) : t \in \mathbb{Z}_+\}$  is recursively defined by (6); and (ii)  $\mathbf{a}(t)$  is independent of  $\{\mathbf{d}(t') : t' > t\}$ , for all  $t \in \mathbb{Z}_+$ . Let  $\Pi$  be the family of all admissible policies.

For  $\mathbf{s} = (\mathbf{x}, \tilde{\mathbf{u}}, R) \in \mathcal{S}$  and  $\mathbf{a} = (\mathbf{y}, \mathbf{z}) \in \mathcal{A}(\mathbf{s})$ , define

$$C(\mathbf{s}, \mathbf{a}) := \mathbf{h}^\top (\mathbf{x} - A^\top \mathbf{y}) + \mathbf{b}^\top (\tilde{\mathbf{u}} - \mathbf{y}). \quad (7)$$

For ease of notation, define the domain of function  $C$  as  $\mathcal{SA} := \{(\mathbf{s}, \mathbf{a}) : \mathbf{s} \in \mathcal{S}, \mathbf{a} \in \mathcal{A}(\mathbf{s})\}$ . The value function, which characterizes the infinite horizon discounted cost associated with policy  $\pi = \{\mathbf{a}(t) : t \in \mathbb{Z}_+\} \in \Pi$  and the initial state  $\mathbf{s}$ , is specified by

$$V^\pi(\mathbf{s}; F) := \mathbb{E} \left[ \sum_{t \in \mathbb{Z}_+} \gamma^t C(\mathbf{s}(t), \mathbf{a}(t)) \mid \mathbf{s}(0) = \mathbf{s} \right], \forall \mathbf{s} \in \mathcal{S}. \quad (8)$$

Here  $\{\mathbf{s}(t)\}$  is the corresponding controlled state process and the expectation is taken with respect to  $\{\mathbf{d}(t) : t \in \mathbb{Z}_+\}$ . In the expression (8), we specify the dependency on the demand distribution  $F$ , as it is the only system primitive that is unknown to the manufacturer. Therefore, an ATO control policy  $\pi$  is optimal if  $V^\pi(\mathbf{s}; F)$  is minimized over  $\pi \in \Pi$ . Define the optimal value function as

$$V^*(\mathbf{s}; F) := \inf_{\pi \in \Pi} V^\pi(\mathbf{s}; F), \forall \mathbf{s} \in \mathcal{S}. \quad (9)$$

Our goal in this paper is to find a *data-driven policy* that approximately solves the data-driven ATO problem. More formally, we aim to find a policy  $\hat{\pi}(\hat{D})$  that can be characterized by i.i.d. sample  $\hat{D}$  from distribution  $F$ , such that the optimality gap

$$V^{\hat{\pi}(\hat{D})}(\mathbf{s}; F) - V^*(\mathbf{s}; F)$$

is small with high probability. This problem is challenging as the optimal value function  $V^*(\cdot; F)$  is multidimensional, has noncompact support, and is also unknown to the manufacturer.

## B Asymmetric Lipschitz Continuity

**Definition 1** (Asymmetric Seminorm). Let  $\rho : \mathbb{R}^p \rightarrow \mathbb{R}_+$  be an asymmetric seminorm if

- (i) For  $\mathbf{x} \in \mathbb{R}^p$  and scalar  $a \geq 0$ ,  $\rho(a\mathbf{x}) = a\rho(\mathbf{x})$ .
- (ii) For  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ ,  $\rho(\mathbf{x} + \mathbf{y}) \leq \rho(\mathbf{x}) + \rho(\mathbf{y})$ .

Asymmetric seminorm is a relaxation of a norm such that  $\rho(-\mathbf{x})$  may not equal  $\rho(\mathbf{x})$  and  $\rho(\mathbf{x}) = 0$  may hold for some  $\mathbf{x} \neq \mathbf{0}$ .

**Definition 2** (Asymmetric Lipschitz Continuity (ALC)). A function  $g : E \subseteq \mathbb{R}^p \rightarrow \mathbb{R}$  is  $L$ -ALC with respect to an asymmetric seminorm  $\rho$  if for  $x_1, x_2 \in E$ ,

$$-L\rho(x_1 - x_2) \leq w(x_2) - w(x_1) \leq L\rho(x_2 - x_1).$$

In general, ALC functions are Lipschitz continuous with possibly different Lipschitz constants, while the opposite is not necessarily true.

For the ATO control problem, consider the following asymmetric seminorms. First, for all  $p \in \mathbb{N}$ , we take asymmetric seminorm  $\rho_0$  on  $\mathbb{R}^p$ , such that  $\rho_0(\mathbf{u}) := (\bar{a} + 1)u_i^- + u_i^+$  for  $\mathbf{u} = (u_1, u_2, \dots, u_p)$ . In particular, we can write

$$\rho_0(\mathbf{s}) = \rho_0(\mathbf{x}) + \rho_0(\tilde{\mathbf{u}}) + \rho_0(R) = \sum_{i=1}^n \rho_0(x_i) + \sum_{j=1}^m \rho_0(\tilde{u}_j) + \sum_{\ell=1}^{L_n} \sum_{i=1}^n \rho_0(R_{i\ell}). \quad (10)$$

We also take a different asymmetric seminorm  $\rho_1$  on  $\mathbb{R}^{p_S + p_A}$ . Let

$$\rho_1(\mathbf{s}, \mathbf{a}) = \rho_0(\mathbf{x} - A\mathbf{y}) + \rho_0(\tilde{\mathbf{u}} - \mathbf{y}) + \rho_0(R) + \rho_0(\mathbf{z})$$

for  $\mathbf{s} = (\mathbf{x}, \tilde{\mathbf{u}}, R)$  and  $\mathbf{a} = (\mathbf{y}, \mathbf{z})$ . It is clear that the functions  $\rho_0$  and  $\rho_1$  are asymmetric seminorms on  $\mathbb{R}^{p_S}$  and  $\mathbb{R}^{p_S + p_A}$ , respectively.

**Theorem 4** (THE ALC PROPERTY). If  $w$  is an  $\bar{c}/(1 - \gamma)$ -ALC function on  $\mathcal{SA}_\alpha$ , then so is  $\mathcal{T}^*w$ .

The goal of the ALC property is to imply Lipschitz continuity of the optimal value function.

**Corollary 1.** For any distribution  $G$  on  $\mathcal{D}$ ,  $V_\alpha^*(\cdot; G)$  is  $L_v$ -Lipschitz continuous on  $\mathcal{S}_\alpha$  for

$$L_v := \frac{\bar{c}(\bar{a} + 1)p_S}{1 - \gamma}.$$

The use of ALC in Theorem 4 is unusual but necessary, as illustrated by the following example.

**Example 1.** Consider an ATO system, where a single product is assembled from two components with unit lead times. Assembly takes one unit of each component. There are no arriving demand. Holding and backlog costs have rate 1. We consider discount factor  $\gamma = 0.9$  and damp rate  $\alpha = 0.1$ .

First, note that Theorem 4 fails under common notions of Lipschitz continuity. Take

$$\begin{aligned} \mathbf{s}^1 &:= (\mathbf{x}^1, \tilde{\mathbf{u}}^1, R^1) = ((1, 1), 1, (0, 0)), & \mathbf{s}^2 &:= (\mathbf{x}^2, \tilde{\mathbf{u}}^2, R^2) = ((2, 2), 2, (0, 0)), \\ \mathbf{a}^1 &:= (\mathbf{y}, \mathbf{z}) = (1, (0, 0)), & \mathbf{a}^2 &:= (\mathbf{y}, \mathbf{z}) = (0, (0, 0)), \quad \text{and} \quad w(\mathbf{s}, \mathbf{a}) := L(x_1 + x_2 + \tilde{u} + y). \end{aligned}$$

Note that the Lipschitz constant depends on the choice of norm. For the commonly used  $L^1$  norm and  $L^\infty$  norm, it is easy to check that  $w$  is  $(L, \|\cdot\|_1)$ -Lipschitz and  $(4L, \|\cdot\|_\infty)$ -Lipschitz. However,

$$\begin{aligned} \|(\mathbf{s}^1, \mathbf{a}^1) - (\mathbf{s}^2, \mathbf{a}^2)\|_\infty &= 1, & \|(\mathbf{s}^1, \mathbf{a}^1) - (\mathbf{s}^2, \mathbf{a}^2)\|_1 &= 4, \\ \mathcal{T}^*w(\mathbf{s}^1, \mathbf{a}^1) &= 0, & \text{and} & \quad T^*w(\mathbf{s}^2, \mathbf{a}^2) = 6 + 4.86L. \end{aligned}$$

There exists no  $L$  such that  $\mathcal{T}^*w$  is  $(L, \|\cdot\|_1)$ -Lipschitz or  $(4L, \|\cdot\|_\infty)$ -Lipschitz continuous.

It is therefore significant that the Lipschitz constant does not necessarily exhibit any contraction property under the Bellman operator  $\mathcal{T}^*$ , despite the discount of cost and the damping of state. In contrast, the function  $w$  can be excluded from Lipschitz continuous classes by considering seminorms. In Theorem 4,  $w$  is not continuous with the proposed asymmetric seminorm  $\rho_1$ .

Next, a seminorm also needs to be asymmetric in order that the Lipschitz constant be preserved by Bellman operator. Take

$$\begin{aligned} \mathbf{s}^3 &= (\mathbf{x}^3, \tilde{\mathbf{u}}^3, R^3) = ((1, 2), 2, (0, 0)), & w_2(\mathbf{s}, \mathbf{a}) &= L(x_2 - x_1 + \tilde{u} - y), \quad \text{and} \\ \rho'_1(\mathbf{s}, \mathbf{a}) &= \|\mathbf{x} - A\mathbf{y}\|_1 + \|\tilde{\mathbf{u}} - \mathbf{y}\|_1 + \sum_{i=1}^n \|R_i\|_1 = |x_1 - y| + |x_2 - y| + |\tilde{u} - y| + |R_1| + |R_2|. \end{aligned}$$

Here seminorm  $\rho'_1$  is the symmetric version of  $\rho_1$ . We can check that

$$|w_2(\mathbf{s}, \mathbf{a})| \leq L\rho'_1(\mathbf{s}, \mathbf{a}), \forall (\mathbf{s}, \mathbf{a}) \in \mathcal{SA}_\alpha,$$

indicating that  $w_2$  is  $(L, \rho'_1)$ -Lipschitz. However,

$$\rho'_1((\mathbf{s}^2, \mathbf{a}^2) - (\mathbf{s}^3, \mathbf{a}^2)) = 1, \quad \mathcal{T}^*w_2(\mathbf{s}_2, \mathbf{a}_2) = 2, \quad \text{and} \quad \mathcal{T}^*w_2(\mathbf{s}_3, \mathbf{a}_2) = 3 + 1.62L.$$

Therefore, there is again no  $L$  such that  $\mathcal{T}^*w_2$  is  $(L, \rho'_1)$ -Lipschitz continuous. In contrast, an  $L$ -ALC function with respect to asymmetric seminorm  $\rho_1$  can decrease by  $2L$  when  $x_1$  increases by 1 (but can increase by at most  $L$ ). We can check that  $w_2$  is  $L$ -ALC. By Theorem 4 that  $\mathcal{T}^*w_2$  is also  $L$ -ALC for sufficiently large  $L$ . Note that the choice of  $\rho_1$  aligns well with the asymmetry between overage and underage in the assembly. An overage of component 1 results in a unit cost from holding the component in the inventory, while an underage of component 1 costs twice as much, i.e., by holding component 2 in the inventory and backlogging one unit of demand. It turns out that, with such asymmetry captured, a sufficiently large ALC constant is preserved by the Bellman's operator.

## C Theoretical Results

**Proposition 1** (BOUNDEDNESS). *For  $\alpha \in (0, 1)$ , there exist proper  $x_\alpha, u_\alpha, R_\alpha, y_\alpha, z_\alpha \in \mathbb{R}_+$  such that  $f_\alpha(\mathbf{s}, \mathbf{a}, \mathbf{d}) \in S_\alpha$  holds whenever  $(\mathbf{s}, \mathbf{a}) \in \mathcal{SA}_\alpha$  and  $\mathbf{d} \in \mathcal{D}$ , where*

$$\begin{aligned} S_\alpha &:= [0, x_\alpha]^n \times [0, u_\alpha]^m \times [0, R_\alpha]^{n \times L_n}, \\ \mathcal{SA}_\alpha &:= \left\{ (\mathbf{s}, \mathbf{a}) : \mathbf{s} \in S_\alpha, \mathbf{a} \in [0, y_\alpha]^m \times [0, z_\alpha]^n, \mathbf{a} \in \mathcal{A}(\mathbf{s}) \right\}. \end{aligned}$$

Moreover,  $(x_\alpha, u_\alpha, R_\alpha, y_\alpha, z_\alpha)$  increases as  $\alpha$  decreases, becomes boundless as  $\alpha \rightarrow 0^+$ .

**Lemma 1.** *For any  $\alpha \in (0, 1)$ ,  $f_\alpha$  is a linear function and  $\mathcal{SA}_\alpha$  is a bounded convex set.*

**Lemma 2** (THE BELLMAN OPERATOR). *For any  $G$  on  $\mathcal{D}$  and  $\alpha \in (0, 1)$ ,*

(i) *the operator  $\mathcal{T}^*$  has a unique fixed point  $Q^*$ , namely,  $\mathcal{T}^*Q^* = Q^*$ .*

(ii) *the function  $V_\alpha^*(\mathbf{s}; G)$  satisfies  $V_\alpha^*(\mathbf{s}; G) = \min_{\mathbf{a} \in \mathcal{A}_\alpha(\mathbf{s})} Q^*(\mathbf{s}, \mathbf{a})$ , for  $\mathbf{s} \in S_\alpha$ .*

**Theorem 5** (PRESERVATION PROPERTY).  *$\mathcal{T}^*w$  is convex whenever  $w : \mathcal{SA}_\alpha \rightarrow \mathbb{R}$  is convex. Moreover,  $\mathcal{T}^*(w) \in \mathcal{C}$  if  $w \in \mathcal{C}$ .*

**Corollary 2.**  *$Q^* \in \mathcal{C}$  and  $V_\alpha^*(\cdot; \widehat{F})$  is convex on  $S_\alpha$ .*

**Proposition 2.** *For any sample  $\widehat{D}$ , Algorithm 1 is well-posed in the sense that*

(i) *there exists a selection  $\xi_K$  that is measurable;*

(ii) *the output policy  $\hat{\pi}^*(\widehat{D}) \in \Pi_\alpha \subseteq \Pi$  is admissible.*

**Lemma 3** (REPRESENTATION POWER OF ICNN). *There exists an ICNN family  $\{Q_\phi : \phi \in \Phi\}$  such that*

$$\inf_{\phi \in \Phi} \|Q - Q_\phi\|_\infty \leq \epsilon$$

*holds for any Lipschitz convex function  $Q \in \mathcal{C}$ .*

## D Proofs

### D.1 Preliminary

In this section, we prove Proposition 1 and Lemma 1.

*Proof.* Proof of Proposition 1. We take

$$u_\alpha = y_\alpha := 2\alpha^{-1}\bar{d}, \quad R_\alpha = z_\alpha := 3\alpha^{-1}\bar{a}L_n\bar{d}, \quad \text{and} \quad x_\alpha := 3\alpha^{-2}\bar{a}L_n\bar{d}.$$

Take an arbitrary  $\mathbf{d} \in [0, \bar{d}]^m$  and  $(\mathbf{s}, \mathbf{a}) \in \mathcal{SA}$  with  $\mathbf{s} = (\mathbf{x}, \tilde{\mathbf{u}}, R)$ ,  $\mathbf{a} = (\mathbf{y}, \mathbf{z})$ . The choice  $(\mathbf{s}, \mathbf{a}) \in \mathcal{SA}$  implies that

$$\begin{aligned} \mathbf{x} &\in [0, x_\alpha]^n, \quad \tilde{\mathbf{u}} \in [0, u_\alpha]^m, \quad R \in [0, R_\alpha]^{n \times L_n}, \\ \mathbf{y} &\in [0, u_\alpha]^n, \quad \mathbf{z} \in [0, z_\alpha]^m \quad \mathbf{y} \leq \tilde{\mathbf{u}}, \quad \text{and} \quad A^\top \mathbf{y} \leq \mathbf{x}. \end{aligned}$$

It holds that

$$\begin{aligned} \mathbf{s}' &= (\mathbf{x}', \tilde{\mathbf{u}}', R') = f_\alpha(\mathbf{s}, \mathbf{a}, \mathbf{d}) \quad \text{for} \\ \mathbf{x}' &= (\mathbf{x} - A^\top \mathbf{y})(1 - \alpha) + R\mathbf{e}_1, \quad \tilde{\mathbf{u}}' = (\tilde{\mathbf{u}} - \mathbf{y})(1 - \alpha) + \mathbf{d}, \quad \text{and} \quad R' = \Gamma(R, \mathbf{z}). \end{aligned} \tag{11}$$

Our goal is to show that

$$\mathbf{x}' \in [0, x_\alpha]^n, \quad \tilde{\mathbf{u}}' \in [0, u_\alpha]^m, \quad \text{and} \quad R' \in [0, R_\alpha]^{n \times L_n}.$$

Plugging the minimal and maximal values of  $(\mathbf{x}, \tilde{\mathbf{u}}, \mathbf{y}, R)$ , it is straightforward to check that

$$\begin{aligned} \mathbf{x}' &\leq (x_\alpha(1 - \alpha) + R_\alpha)\mathbf{1} = x_\alpha\mathbf{1}, \quad \tilde{\mathbf{u}}' \leq (u_\alpha(1 - \alpha) + \bar{d})\mathbf{1} \leq u_\alpha\mathbf{1}, \\ \mathbf{x}' &= (\mathbf{x} - A^\top \mathbf{y})(1 - \alpha) \geq \mathbf{0}, \quad \text{and} \quad \tilde{\mathbf{u}}' = (\tilde{\mathbf{u}} - \mathbf{y})\alpha \geq \mathbf{0}, \end{aligned}$$

where the last two inequalities utilize  $A^\top \mathbf{y} \leq \mathbf{x}$  and  $\mathbf{y} \leq \tilde{\mathbf{u}}$ , respectively. As  $\mathbf{z}$ ,  $R$ ,  $U_1$ , and  $U_2$  have nonnegative entries, it is clear that

$$R' = RU_1 + \text{diag}(\mathbf{z})U_2 \geq \mathbf{0}.$$

Finally, the choices of  $U_1$ ,  $U_2$ , and  $R \in \mathcal{R}$  imply that  $RU_1$  and  $\text{diag}(\mathbf{z})U_2$  has no common nonzero entry, so that  $R \in [0, R_\alpha]^{n \times (L_n - 1)}$  and  $z_\alpha = R_\alpha$  imply

$$(R')_{j,\ell} = (RU_1 + \text{diag}(\mathbf{z})U_2)_{j,\ell} \leq R_\alpha, \quad \text{for all } j \in [n] \quad \text{and} \quad \ell \in [L_n].$$

Therefore,  $R' \in [0, R_\alpha]^{n \times L_n}$ , finishing the proof.  $\square$

*Proof of Lemma 1.* Let  $U_1$  be an  $L_n \times L_n$  0-1 matrix such that  $(U_1)_{ij} = 1$  if  $j = i - 1$ . Let  $U_2$  be an  $n \times L_n$  0-1 matrix such that  $(U_2)_{ij} = 1$  when  $j = L_i - 1$ . Define linear mapping

$$\Gamma(R, \mathbf{z}) := RU_1 + \text{diag}(\mathbf{z})U_2.$$

Here,  $\text{diag}(\mathbf{z})$  returns an  $n$  by  $n$  diagonal matrix generated from  $\mathbf{z}$ . We can write

$$\begin{aligned} f(\mathbf{s}, \mathbf{a}, \mathbf{d}) &:= \left( (\mathbf{x} - A^\top \mathbf{y}) + R\mathbf{e}_1, (\tilde{\mathbf{u}} - \mathbf{y}) + \mathbf{d}, \Gamma(R, \mathbf{z}) \right), \quad \text{for} \\ \mathbf{s} &= (\mathbf{x}, \tilde{\mathbf{u}}, R) \in \mathbb{R}^{ps}, \quad \mathbf{a} = (\mathbf{y}, \mathbf{z}) \in \mathbb{R}^{pA}, \quad \text{and} \quad \mathbf{d} \in \mathbb{R}^m. \end{aligned} \tag{12}$$

In particular, mapping  $f$  is linear.

Next, the boundedness of  $\mathcal{SA}_\alpha$  clearly follows from Proposition 1. It follows from

$$\mathcal{SA}_\alpha = \{(\mathbf{s}, \mathbf{a}) \in \mathcal{S}_\alpha \times \mathcal{A}_\alpha : A^\top \mathbf{y} \leq \mathbf{x} \quad \text{and} \quad \mathbf{y} \leq \tilde{\mathbf{u}}\}$$

that  $\mathcal{SA}_\alpha$  is a convex polyhedron.  $\square$

### D.2 The Shape Preservation Properties

In this section, we present our original proofs to Theorems 5 and 4, which establish that Bellman operator  $\mathcal{T}^*$  preserves boundedness, convexity, and asymmetric Lipschitz continuity.

*Proof.* Proof of Theorem 5. We first show the boundedness of the operator  $\mathcal{T}^*$ . Take

$$B_\alpha = \bar{c}(m + n)x_\alpha / (1 - \gamma).$$

For  $(\mathbf{s}, \mathbf{a}) \in \mathcal{SA}_\alpha$ , (7) gives

$$0 \leq C(\mathbf{s}, \mathbf{a}) = \mathbf{h}^\top (\mathbf{x} - A^\top \mathbf{y}) + \mathbf{b}^\top (\tilde{\mathbf{u}} - \mathbf{y}_1) \leq m x_\alpha \|\mathbf{h}\|_\infty + n u_\alpha \|\mathbf{b}\|_\infty, \quad \text{for all } (\mathbf{s}, \mathbf{a}) \in \mathcal{SA}_\alpha.$$

$$(\mathcal{T}^* w)(\mathbf{s}, \mathbf{a}) := C(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\hat{\mathbf{d}} \sim G} \left[ \min_{\mathbf{a}' \in \mathcal{A}_\alpha(\mathbf{s}')} w(\mathbf{s}', \mathbf{a}') \right], \quad \text{for } \mathbf{s}' = f_\alpha(\mathbf{s}, \mathbf{a}, \hat{\mathbf{d}}), \quad \mathbf{s} \in \mathcal{S}_\alpha. \quad (13)$$

Together with (13) and  $u_\alpha \leq x_\alpha$ ,

$$\mathcal{T}^*[w(\mathbf{s}, \mathbf{a})] = C(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{d}} \left[ \min_{\mathbf{a}' \in \mathcal{A}_\alpha(\mathbf{s}')} w(\mathbf{s}', \mathbf{a}') \right] \leq \bar{c}(m+n)x_\alpha + \gamma B_\alpha \leq B_\alpha, \quad \text{for all } (\mathbf{s}, \mathbf{a}) \in \mathcal{SA}_\alpha.$$

Therefore,  $\mathcal{T}^* w$  is bounded by  $B_\alpha$  if  $w$  is bounded by  $B_\alpha$ .

For the convex property of  $\mathcal{T}^*$ , we take an arbitrary convex function  $w$  on  $\mathcal{SA}_\alpha$ , take arbitrary state and action pairs  $(\mathbf{s}_i, \mathbf{a}_i) \in \mathcal{SA}_\alpha$ ,  $i = 1, 2$  and let

$$\mathbf{s}_3 = \frac{\mathbf{s}_1 + \mathbf{s}_2}{2} \quad \text{and} \quad \mathbf{a}_3 = \frac{\mathbf{a}_1 + \mathbf{a}_2}{2}.$$

Temporarily write

$$\mathbf{s}'_i := f(\mathbf{s}_i, \mathbf{a}_i, \mathbf{d}), \quad c_i := C(\mathbf{s}_i, \mathbf{a}_i), \quad \text{for } i = 1, 2, 3, \quad \mathbf{a}'_i \in \arg \min_{\mathbf{a} \in \mathcal{A}_\alpha(\mathbf{s}'_i)} w(\mathbf{s}'_i, \mathbf{a}) \quad \text{for } i = 1, 2.$$

Note that the minimal is obtained as  $w(\mathbf{s}'_i, \cdot)$  is convex on the compact set  $\mathcal{A}_\alpha(\mathbf{s}'_i)$ . It follows from the linearity of  $f$  and  $C$  that

$$\mathbf{s}'_3 = (\mathbf{s}'_1 + \mathbf{s}'_2)/2 \quad \text{and} \quad c_3 = (c_1 + c_2)/2.$$

Clearly  $(\mathbf{s}'_i, \mathbf{a}'_i) \in \mathcal{SA}_\alpha$  for  $i = 1, 2$ . It follows from Lemma 1 that  $\mathcal{SA}_\alpha$  is convex, implying that  $(\mathbf{s}'_3, \mathbf{a}'_3) \in \mathcal{SA}_\alpha$  and, in turn,  $\mathbf{a}'_3 \in \mathcal{A}_\alpha(\mathbf{s}'_3)$  for  $\mathbf{a}'_3 := (\mathbf{a}'_1 + \mathbf{a}'_2)/2$ . Therefore,

$$\begin{aligned} \mathcal{T}^* w(\mathbf{s}_1, \mathbf{a}_1) + \mathcal{T}^* w(\mathbf{s}_2, \mathbf{a}_2) &= c_1 + c_2 + \gamma \mathbb{E}_{\mathbf{d}} [w(\mathbf{s}'_1, \mathbf{a}'_1) + w(\mathbf{s}'_2, \mathbf{a}'_2)] \\ &\geq c_1 + c_2 + 2\gamma \mathbb{E}_{\mathbf{d}} [w(\mathbf{s}'_3, \mathbf{a}'_3)] \geq 2c_3 + 2\gamma \mathbb{E}_{\mathbf{d}} \left[ \min_{\mathbf{a} \in \mathcal{A}_\alpha(\mathbf{s}'_3)} w(\mathbf{s}'_3, \mathbf{a}) \right] = 2\mathcal{T}^* w(\mathbf{s}_3, \mathbf{a}_3), \end{aligned}$$

in which first inequality results from the convexity of  $w$ . Because the choice of  $(\mathbf{s}_i, \mathbf{a}_i)$  is arbitrary,  $\mathcal{T}^* w$  is convex on  $\mathcal{SA}_\alpha$  whenever  $w$  is convex on  $\mathcal{SA}_\alpha$ , finishing the proof of Theorem 5.  $\square$

The proof of Theorem 4 is based on Lemma 4, a supporting result that relates the asymmetric norms  $\rho_1$  to  $\rho_0$ .

**Lemma 4.** For any  $\alpha \in [0, 1]$ ,  $\mathbf{s}_1 \in \mathcal{S}_\alpha$ ,  $\mathbf{s}_2 \in \mathcal{S}_\alpha$ , it holds that

$$\max_{\mathbf{a}_1 \in \mathcal{A}_\alpha(\mathbf{s}_1)} \min_{\mathbf{a}_2 \in \mathcal{A}_\alpha(\mathbf{s}_2)} \rho_1((\mathbf{s}_2 - \mathbf{s}_1, \mathbf{a}_2 - \mathbf{a}_1)) \leq \rho_0(\mathbf{s}_2 - \mathbf{s}_1). \quad (14)$$

Note that the maximum and minimum is attained as  $\rho_1$  is a continuous function and  $\mathcal{A}_\alpha(\mathbf{s})$  is compact for  $\mathbf{s} \in \mathcal{S}_\alpha$ .

*Proof.* Proof of Theorem 4. Take an arbitrary  $L$ -ALC function  $w$  on  $\mathcal{SA}_\alpha$ ,  $(\mathbf{s}_1, \mathbf{a}_1) \in \mathcal{SA}_\alpha$ , and  $(\mathbf{s}_2, \mathbf{a}_2) \in \mathcal{SA}_\alpha$ .

First, for the single period cost function  $C$  in (7), we can compute that.

$$\begin{aligned} C(\mathbf{s}_2, \mathbf{a}_2) - C(\mathbf{s}_1, \mathbf{a}_1) &= \mathbf{h}^\top (\mathbf{x}_2 - \mathbf{x}_1) - \mathbf{h}^\top A^\top (\mathbf{y}_2 - \mathbf{y}_1) + \mathbf{b}^\top (\tilde{\mathbf{u}}_2 - \tilde{\mathbf{u}}_1 - \mathbf{y}_2 + \mathbf{y}_1) \\ &\leq \bar{c} \rho_0(\mathbf{x}_2 - \mathbf{x}_1 - A^\top (\mathbf{y}_2 - \mathbf{y}_1)) + \bar{c} \rho_0(\tilde{\mathbf{u}}_2 - \tilde{\mathbf{u}}_1 - \mathbf{y}_2 + \mathbf{y}_1) \\ &\leq \bar{c} \rho_1(\mathbf{s}_2 - \mathbf{s}_1, \mathbf{a}_2 - \mathbf{a}_1). \end{aligned}$$

Second, we temporarily denote  $\mathbf{s}'_i := f_\alpha(\mathbf{s}_i, \mathbf{a}_i, \mathbf{d})$  for  $i = 1, 2$ . By (13),

$$\begin{aligned} &\mathcal{T}^* w(\mathbf{s}_2, \mathbf{a}_2) - \mathcal{T}^* w(\mathbf{s}_1, \mathbf{a}_1) \\ &= C(\mathbf{s}_2, \mathbf{a}_2) - C(\mathbf{s}_1, \mathbf{a}_1) + \gamma \mathbb{E} \left[ \min_{\mathbf{a} \in \mathcal{A}_\alpha(\mathbf{s}'_2)} w(\mathbf{s}'_2, \mathbf{a}) - \min_{\mathbf{a} \in \mathcal{A}_\alpha(\mathbf{s}'_1)} w(\mathbf{s}'_1, \mathbf{a}) \right] \\ &= C(\mathbf{s}_2, \mathbf{a}_2) - C(\mathbf{s}_1, \mathbf{a}_1) + \gamma \mathbb{E} \left[ \max_{\mathbf{a}'_1 \in \mathcal{A}_\alpha(\mathbf{s}'_1)} \min_{\mathbf{a}'_2 \in \mathcal{A}_\alpha(\mathbf{s}'_2)} w(\mathbf{s}'_2, \mathbf{a}'_2) - w(\mathbf{s}'_1, \mathbf{a}'_1) \right] \\ &\leq \bar{c} \rho_1(\mathbf{s}_2 - \mathbf{s}_1, \mathbf{a}_2 - \mathbf{a}_1) + \gamma L \mathbb{E} \left[ \max_{\mathbf{a}'_1 \in \mathcal{A}_\alpha(\mathbf{s}'_1)} \min_{\mathbf{a}'_2 \in \mathcal{A}_\alpha(\mathbf{s}'_2)} \rho_1(\mathbf{s}'_2 - \mathbf{s}'_1, \mathbf{a}'_2 - \mathbf{a}'_1) \right] \\ &\leq \bar{c} \rho_1(\mathbf{s}_2 - \mathbf{s}_1, \mathbf{a}_2 - \mathbf{a}_1) + \gamma L \mathbb{E} [\rho_0(\mathbf{s}'_2 - \mathbf{s}'_1)] \\ &= \bar{c} \rho_1(\mathbf{s}_2 - \mathbf{s}_1, \mathbf{a}_2 - \mathbf{a}_1) + \gamma L \rho_0(f_\alpha(\mathbf{s}_2 - \mathbf{s}_1, \mathbf{a}_2 - \mathbf{a}_1, \mathbf{0})). \end{aligned} \quad (15)$$

Note that the first inequality follows from the fact that  $C$  is  $\bar{c}$ -ALC and  $w$  is  $L$ , the second inequality results from Lemma 4, and the last equality in (15) uses the linearity of  $f_\alpha$ .

Next, we expand the expression of  $f_\alpha$  in the right-hand side of (15). For

$$(\mathbf{x}', \tilde{\mathbf{u}}', R') := f_\alpha(\mathbf{s}_2 - \mathbf{s}_1, \mathbf{a}_2 - \mathbf{a}_1, \mathbf{0}),$$

the expressions in (11) gives

$$\begin{aligned}\mathbf{x}' &= (1 - \alpha)(\mathbf{x}_2 - \mathbf{x}_1 - A^\top(\mathbf{y}_2 - \mathbf{y}_1)) + (R_2 - R_1)\mathbf{e}_1, \\ \tilde{\mathbf{u}}' &= (1 - \alpha)(\tilde{\mathbf{u}}_2 - \tilde{\mathbf{u}}_1 - (\mathbf{y}_2 - \mathbf{y}_1)), \quad \text{and} \\ R' &= (R_2 - R_1)U_1 + U_2 \text{diag}(\mathbf{z}_2 - \mathbf{z}_1).\end{aligned}$$

Using the triangular inequality for asymmetric seminorm  $\rho_0$  and dropping the scalar  $(1 - \alpha)$ ,

$$\begin{aligned}\rho_0(\mathbf{x}') &\leq \rho_0(\mathbf{x}_2 - \mathbf{x}_1 - A^\top(\mathbf{y}_2 - \mathbf{y}_1)) + \rho_0((R_2 - R_1)\mathbf{e}_1), \\ \rho_0(\tilde{\mathbf{u}}') &\leq \rho_0(\tilde{\mathbf{u}}_2 - \tilde{\mathbf{u}}_1 - (\mathbf{y}_2 - \mathbf{y}_1)), \quad \text{and} \\ \rho_0(R') &\leq \rho_0((R_2 - R_1)U_1) + \rho_0(U_2 \text{diag}(\mathbf{z}_2 - \mathbf{z}_1)).\end{aligned}$$

Expanding the expressions of matrices  $U_1$  and  $U_2$ ,

$$\rho_0((R_2 - R_1)U_1) = \sum_{\ell=2}^{L_n} \rho_0((R_2 - R_1)\mathbf{e}_\ell) = \rho_0(R_2 - R_1) \quad \text{and} \quad \rho_0(U_2 \text{diag}(\mathbf{z}_2 - \mathbf{z}_1)) = \rho_0(\mathbf{z}_2 - \mathbf{z}_1).$$

To summarize, we have shown that

$$\begin{aligned}&\rho_0(f_\alpha(\mathbf{s}_2 - \mathbf{s}_1, \mathbf{a}_2 - \mathbf{a}_1, \mathbf{0})) \\ &= \rho_0(\mathbf{x}') + \rho_0(\tilde{\mathbf{u}}') + \rho_0(R') \\ &\leq \rho_0(\mathbf{x}_2 - \mathbf{x}_1 - A^\top(\mathbf{y}_2 - \mathbf{y}_1)) + \rho_0(\tilde{\mathbf{u}}_2 - \tilde{\mathbf{u}}_1 - (\mathbf{y}_2 - \mathbf{y}_1)) + \rho_0(\mathbf{z}_2 - \mathbf{z}_1) + \rho_0(R_2 - R_1) \\ &= \rho_1(\mathbf{s}_2 - \mathbf{s}_1, \mathbf{a}_2 - \mathbf{a}_1).\end{aligned}$$

Together with (15),

$$\mathcal{T}^*w(\mathbf{s}_2, \mathbf{a}_2) - \mathcal{T}^*w(\mathbf{s}_1, \mathbf{a}_1) \leq (\bar{c} + \gamma L)\rho_1(\mathbf{s}_1 - \mathbf{s}_2, \mathbf{a}_1 - \mathbf{a}_2), \quad \text{for all } (\mathbf{s}_i, \mathbf{a}_i) \in \mathcal{S}\mathcal{A}_\alpha, \quad i = 1, 2.$$

implying that  $\mathcal{T}w$  is a  $(\bar{c} + \gamma L)$ -ALC function.  $\square$

The proof of Lemma 4 is based on the following linear algebra result.

**Lemma 5.** Let  $p$  by  $q$  matrix  $H = (h_{ij})$  satisfy  $h_{ij} \in \{0\} \cup [1, \infty)$  for all  $i \in [p]$  and  $j \in [q]$ . Then for all  $\mathbf{w} \in \mathbb{R}_+^q$  and  $\mathbf{b} \in \mathbb{R}_+^p$  such that  $\mathbf{b} \leq H\mathbf{w}$ , there exists  $0 \leq \mathbf{v} \leq \mathbf{w}$  such that

$$H\mathbf{v} \leq \mathbf{b} \quad \text{and} \quad \|\mathbf{b} - H\mathbf{v}\|_1 \leq \|\mathbf{1}^\top H\|_\infty \|\mathbf{H}\mathbf{w} - \mathbf{b}\|_1. \quad (16)$$

*Proof of Lemma 4.* Fix states  $\mathbf{s}_1 := (\mathbf{x}_1, \tilde{\mathbf{u}}_1, R_1) \in \mathcal{S}_\alpha$ ,  $\mathbf{s}_2 := (\mathbf{x}_2, \tilde{\mathbf{u}}_2, R_2) \in \mathcal{S}_\alpha$ , an arbitrary action  $\mathbf{a}_1 := (\mathbf{y}_1, \mathbf{z}_1) \in \mathcal{A}_\alpha(\mathbf{s}_1)$ . Our goal is to find an action  $\mathbf{a}_2 = (\mathbf{y}_2, \mathbf{z}_2) \in \mathcal{A}_\alpha(\mathbf{s}_2)$  such that

$$\rho_1(\mathbf{s}_2 - \mathbf{s}_1, \mathbf{a}_2 - \mathbf{a}_1) \leq \rho_0(\mathbf{s}_2 - \mathbf{s}_1). \quad (17)$$

We take  $\mathbf{z}_2 = \mathbf{z}_1$  and find  $\mathbf{y}_2$  by Lemma 5. Let  $\mathbf{b}_i := (\mathbf{x}_i, \tilde{\mathbf{u}}_i)$  for  $i = 1, 2$ ,  $H = [A, I]^\top$  and  $\mathbf{b}'_2 = \min\{H\mathbf{y}_1, \mathbf{b}_2\}$ . It is straightforward to check that  $H$  satisfies the condition for Lemma 5. It is also clear that  $\mathbf{b}'_2 \leq H\mathbf{y}_1 \in \mathbb{R}_+^{m+n}$  and  $\mathbf{y}_2 \in \mathbb{R}_+^m$ . In Lemma 5 we take  $\mathbf{b} = \mathbf{b}'_2$  and  $\mathbf{w} = \mathbf{y}_1$ , then there exists  $0 \leq \mathbf{y}_2 \leq \mathbf{y}_1 \in \mathbb{R}_+^m$  such that

$$H\mathbf{y}_2 \leq \mathbf{b}'_2, \quad \text{and} \quad \mathbf{1}^\top(\mathbf{b}'_2 - H\mathbf{y}_2) \leq \|\mathbf{1}^\top H\|_\infty (H\mathbf{y}_1 - \mathbf{b}'_2). \quad (18)$$

Because  $H\mathbf{y}_2 \leq \mathbf{b}'_2 \leq \mathbf{b}_2$ , it is easy to check that  $A^\top \mathbf{y}_2 \leq \mathbf{x}_2$  and  $\mathbf{y}_2 \leq \tilde{\mathbf{u}}_2$ , implying that  $\mathbf{a}_2 \in \mathcal{A}_\alpha(\mathbf{s}_2)$ . To show (17), the choice of  $\rho_1$  gives

$$\rho_1(\mathbf{s}_2 - \mathbf{s}_1, \mathbf{a}_2 - \mathbf{a}_1) = \rho_0(\mathbf{b}_2 - \mathbf{b}_1 - H(\mathbf{y}_2 - \mathbf{y}_1)) + \rho_0(R_2 - R_1). \quad (19)$$

For the first term on the right-hand side of (19), the triangular inequality of  $\rho_0$  gives

$$\rho_0(\mathbf{b}_2 - \mathbf{b}_1 - H(\mathbf{y}_2 - \mathbf{y}_1)) \leq \rho_0(\mathbf{b}_2 - \mathbf{b}'_2 - \mathbf{b}_1 + H\mathbf{y}_1) + \rho_0(\mathbf{b}'_2 - H\mathbf{y}_2).$$

To obtain an upper bound of  $\rho_0(\mathbf{b}'_2 - H\mathbf{y}_2)$ , the second inequality in (18) gives  $\mathbf{b}'_2 - H\mathbf{y}_2 \geq \mathbf{0}$ . It follows from the choice of  $\rho_0$  and the second inequality in (18) that

$$\rho_0(\mathbf{b}'_2 - H\mathbf{y}_2) = \|\mathbf{b}'_2 - H\mathbf{y}_2\|_1 \leq \|\mathbf{1}^\top H\|_\infty \|H\mathbf{y}_1 - \mathbf{b}'_2\|_1,$$

Because  $\mathbf{b}'_2 - H\mathbf{y}_1 \leq \mathbf{0}$ , the choice of  $\rho_0$  gives

$$\rho_0(\mathbf{b}'_2 - H\mathbf{y}_1) = \|\mathbf{1}^\top H\|_\infty (H\mathbf{y}_1 - \mathbf{b}'_2),$$

and thus

$$\rho_0(\mathbf{b}'_2 - H\mathbf{y}_2) \leq \rho_0(\mathbf{b}'_2 - H\mathbf{y}_1).$$

Therefore, we have shown that

$$\rho_0(\mathbf{b}_2 - \mathbf{b}_1 - H(\mathbf{y}_2 - \mathbf{y}_1)) \leq \rho_0(\mathbf{b}'_2 - H\mathbf{y}_1) + \rho_0(\mathbf{b}_2 - \mathbf{b}'_2 - \mathbf{b}_1 + H\mathbf{y}_1).$$

We next show that the two vectors on the right-hand side have the same sign in every entry, i.e.,

$$(\mathbf{b}'_2 - H\mathbf{y}_1)_\ell (\mathbf{b}_2 - \mathbf{b}'_2 - \mathbf{b}_1 + H\mathbf{y}_1)_\ell \geq 0, \quad \text{for all } \ell \in [m+n], \quad (20)$$

It is clear that  $\mathbf{b}'_2 - H\mathbf{y}_1 = \min\{0, \mathbf{b}_2 - H\mathbf{y}_1\} \leq \mathbf{0}$ . For any index  $\ell \in [m+n]$  such that  $\mathbf{b}'_2 - H\mathbf{y}_1 = \min\{0, \mathbf{b}_2 - H\mathbf{y}_1\}$  has a strictly negative  $\ell$ th entry, it must hold that

$$(\mathbf{b}_2 - H\mathbf{y}_1)_\ell = (\mathbf{b}'_2 - H\mathbf{y}_1)_\ell < 0.$$

As  $H\mathbf{y}_1 \leq \mathbf{b}_1$  by  $\mathbf{a}_1 \in \mathcal{A}_\alpha(\mathbf{s}_1)$ , we have shown that  $(\mathbf{b}_2 - \mathbf{b}'_2 - \mathbf{b}_1 + H\mathbf{y}_1)_\ell = (H\mathbf{y}_1 - \mathbf{b}_1)_\ell \leq 0$ , giving (20).

Now the choice of  $\rho_0$  and (20) imply that

$$\rho_0(\mathbf{b}'_2 - H\mathbf{y}_1) + \rho_0(\mathbf{b}_2 - \mathbf{b}'_2 - \mathbf{b}_1 + H\mathbf{y}_1) = \rho_0(\mathbf{b}_2 - \mathbf{b}_1),$$

and we have thus computed that

$$\rho_0(\mathbf{b}_2 - \mathbf{b}_1 - H(\mathbf{y}_2 - \mathbf{y}_1)) \leq \rho_0(\mathbf{b}_2 - \mathbf{b}_1) = \rho_0(\mathbf{x}_2 - \mathbf{x}_1) + \rho_0(\tilde{\mathbf{u}}_2 - \tilde{\mathbf{u}}_1).$$

Together with (19) and (10), we obtain (17), finishing the proof.  $\square$

*Proof of Lemma 5.* Take  $\mathbf{w}_0 = \mathbf{w}$  and recursively define the sequence  $\{\mathbf{w}_i : i \in [p]\}$  such that  $\mathbf{w}_i$  maximizes  $\mathbf{1}^\top \mathbf{w}'$  over  $0 \leq \mathbf{w}' \leq \mathbf{w}_{i-1}$  subject to  $(H\mathbf{w}')_i \leq b_i$ . Here  $b_i$  is the  $i$ th entry of vector  $\mathbf{b}$ . We claim that  $\mathbf{v} = \mathbf{w}_p$  satisfies (16). To show this, it follows from  $\mathbf{w}_p \leq \mathbf{w}_i$  and that  $H$  is a nonnegative matrix that  $(H\mathbf{v})_i = (H\mathbf{w}_p)_i \leq (H\mathbf{w}_i)_i \leq b_i$  for all  $i \in [p]$ , indicating that  $H\mathbf{v} \leq \mathbf{b}$ .

To prove the second inequality in (16), it follows from  $H\mathbf{v} \leq \mathbf{b} \leq H\mathbf{w}$  that

$$\begin{aligned} \|\mathbf{b} - H\mathbf{v}\|_1 &\leq \|H\mathbf{w} - H\mathbf{v}\|_1 = \mathbf{1}^\top H(\mathbf{w} - \mathbf{v}) \leq \|\mathbf{1}^\top H\|_\infty \|\mathbf{w} - \mathbf{v}\|_1 \\ &= \sum_{i=1}^p \|\mathbf{1}^\top H\|_\infty \|\mathbf{w}_{i-1} - \mathbf{w}_i\|_1. \end{aligned} \quad (21)$$

Here, the last equality follows from  $\mathbf{w} = \mathbf{w}_0 \geq \mathbf{w}_1 \geq \dots \geq \mathbf{w}_p = \mathbf{v}$ . Because  $\mathbf{w}_i$  maximizes  $\mathbf{1}^\top \mathbf{w}'$  subject to  $0 \leq \mathbf{w}' \leq \mathbf{w}_{i-1}$  and  $(H\mathbf{w}')_i = \sum_{j=1}^q h_{ij} w'_j \leq b_i$ , the optimality condition gives

$$1\{h_{ij} = 0\}(\mathbf{w}_i - \mathbf{w}_{i-1})_j = 0 \quad \text{for all } j \in [q] \quad \text{and} \quad (H\mathbf{w}_i)_i = \min\{b_i, (H\mathbf{w}_{i-1})_i\}.$$

Together with the fact that  $h_{ij} \in \{0\} \cup [1, \infty)$  gives

$$\|\mathbf{w}_{i-1} - \mathbf{w}_i\|_1 = \sum_{j \in [q]} 1\{h_{ij} \neq 0\}(\mathbf{w}_{i-1} - \mathbf{w}_i)_j \leq \sum_{j \in [q]} h_{ij}(\mathbf{w}_{i-1} - \mathbf{w}_i) = (H\mathbf{w}_{i-1} - H\mathbf{w}_i)_i.$$

We can compute the right-hand side by

$$(H\mathbf{w}_{i-1} - H\mathbf{w}_i)_i = (H\mathbf{w}_{i-1})_i - (H\mathbf{w}_i)_i = (H\mathbf{w}_{i-1})_i - \min\{(H\mathbf{w}_{i-1})_i, b_i\} = \max\{0, (H\mathbf{w}_{i-1})_i - b_i\}.$$

Together with  $\mathbf{w}_{i-1} \leq \mathbf{w}$  and  $H\mathbf{w} \geq \mathbf{b}$ , we have shown that

$$\|\mathbf{w}_{i-1} - \mathbf{w}_i\|_1 \leq (H\mathbf{w}_{i-1} - H\mathbf{w}_i)_i = \max\{0, (H\mathbf{w}_{i-1})_i - b_i\} \leq (H\mathbf{w} - \mathbf{b})_i, \quad \text{for all } i \in [p].$$

Plugging in (21), we obtain the second inequality in (16), finishing the proof.  $\square$

### D.3 Contraction of Bellman Operators

In this section, we state and prove Proposition 3, an extended version of Lemma 2. The statements in Proposition 3 are mostly standard. The proofs of Corollaries 2 and 1 also appear at this section.

We introduce the following notations.

$$\mathcal{A}_\alpha(\mathbf{s}) := \{\mathbf{a} \in \mathbb{R}^{p_A} : (\mathbf{s}, \mathbf{a}) \in \mathcal{SA}_\alpha\} = \mathcal{A}(\mathbf{s}) \cap ([0, y_\alpha]^m \times [0, z_\alpha]^n) \quad \text{for } \mathbf{s} \in \mathcal{S}_\alpha. \quad (22)$$

First, define Markovian control  $\xi$  by a Borel measurable map from  $\mathcal{S}_\alpha$  to  $\mathbb{R}_+^{p_A}$ , such that  $\xi(\mathbf{s}) \in \mathcal{A}_\alpha(\mathbf{s})$  for  $\mathbf{s} \in \mathcal{S}_\alpha$ , where we recall  $\mathcal{A}_\alpha$  from (22).

$$\mathbf{s}_\alpha(t+1) = f_\alpha(\mathbf{s}_\alpha(t), \mathbf{a}(t), \mathbf{d}(t+1)), \quad \text{for all } t \geq 0. \quad (23)$$

We say that policy  $\pi \in \Pi_\alpha$  is induced by  $\xi$  if  $\pi = \{\mathbf{a}(t)\}$  is recursively characterized by  $\mathbf{a}(t) = \xi(\mathbf{s}_\alpha(t))$  and  $\{\mathbf{s}_\alpha(t)\}$  from (23). Second, define Bellman operator  $\mathcal{T}^\xi$  for Markovian policy  $\xi \in \Xi$  by

$$\mathcal{T}^\xi[w(\mathbf{s}, \mathbf{a})] := C(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{d} \sim G}[w(\mathbf{s}', \xi(\mathbf{s}'))], \quad (24)$$

for  $(\mathbf{s}, \mathbf{a}) \in \mathcal{SA}_\alpha$  and  $\mathbf{s}' = f(\mathbf{s}, \mathbf{a}, \mathbf{d})$ . The following result summarizes all the properties that we obtain for the Bellman operators  $\mathcal{T}^*$  and  $\mathcal{T}^\xi$  for later reference.

**Proposition 3.** For any  $\alpha \in (0, 1)$  and distribution  $G$  on  $\mathcal{D}$ ,

(i) operators  $\mathcal{T}^\xi$  and  $\mathcal{T}^*$  are contraction mappings, i.e., for any  $w_1$  and  $w_2$  that are measurable on  $\mathcal{SA}_\alpha$ ,

$$\|\mathcal{T}^\xi(w_1) - \mathcal{T}^\xi(w_2)\|_\infty \leq \gamma \|w_1 - w_2\|_\infty \quad \text{and} \quad \|\mathcal{T}^*(w_1) - \mathcal{T}^*(w_2)\|_\infty \leq \gamma \|w_1 - w_2\|_\infty.$$

Here, the infinite norm  $\|\cdot\|_\infty$  is defined by

$$\|w\|_\infty = \sup_{(\mathbf{s}, \mathbf{a}) \in \mathcal{SA}_\alpha} |w(\mathbf{s}, \mathbf{a})|;$$

(ii) operator  $\mathcal{T}^\xi$  has a unique fixed point  $Q^\xi : \mathcal{SA}_\alpha \rightarrow \mathbb{R}$ , i.e.,  $\mathcal{T}^\xi(Q^\xi) = Q^\xi$ . Moreover,  $Q^\xi$  satisfies  $V_\alpha^\pi(\mathbf{s}; G) = Q^\xi(\mathbf{s}, \xi(\mathbf{s}))$  for all  $\mathbf{s} \in \mathcal{S}_\alpha$ , where  $\pi$  is the induced policy by  $\xi$ .

(iii) operator  $\mathcal{T}^*$  has a unique fixed point  $Q^* \in \mathcal{C}$ . For  $w_0 \in \mathcal{C}$ , the recursively defined sequence  $\{w_k : k \geq 0\}$  by  $w_{k+1} := \mathcal{T}^* w_k$  satisfies  $w_k \in \mathcal{C}$  and

$$\|w_k - Q^*\|_\infty \leq \gamma^k \|w_0 - Q^*\|_\infty, \quad \text{for } k \in \mathbb{Z}_+. \quad (25)$$

(iv) there exists an optimal Markovian control  $\xi^* \in \Xi$  satisfying  $Q^{\xi^*} = Q^*$  and

$$V_\alpha^{\pi^*}(\mathbf{s}; G) = V_\alpha^*(\mathbf{s}; G) = \min_{\mathbf{a} \in \mathcal{A}_\alpha(\mathbf{s})} Q^*(\mathbf{s}, \mathbf{a}), \quad \text{for all } \mathbf{s} \in \mathcal{S}_\alpha. \quad (26)$$

Here  $\pi^* \in \Pi_\alpha$  is the induced policy by  $\xi^*$ .

Note that all results in Corollary 2 and Lemma 2 immediately follow from Proposition 3.

*Proof.* Proof of Proposition 3.

**Proof of Assertion (i):** It follows from Blackwell's sufficient conditions (see [5]) that assertion (i) is implied by the following monotonicity and discounting properties of the operator  $\mathcal{T} \in \{\mathcal{T}^\xi, \mathcal{T}^*\}$ :

- *Monotonicity:* If  $w(\mathbf{s}, \mathbf{a}) \geq w'(\mathbf{s}, \mathbf{a})$  for all  $(\mathbf{s}, \mathbf{a}) \in \mathcal{SA}_\alpha$  where  $w, w'$  are measurable functions on  $\mathcal{SA}_\alpha$ , then  $\mathcal{T}(w)(\mathbf{s}, \mathbf{a}) \geq \mathcal{T}(w')(\mathbf{s}, \mathbf{a})$ .
- *Discounting:* For all measurable function  $w$  on  $\mathcal{SA}_\alpha$ ,  $\mathcal{T}(w + b)(\mathbf{s}, \mathbf{a}) \leq \mathcal{T}(w)(\mathbf{s}, \mathbf{a}) + \gamma b$  for all  $b \geq 0$  and  $(\mathbf{s}, \mathbf{a}) \in \mathcal{SA}_\alpha$ . Here  $(w + b)(\mathbf{s}, \mathbf{a}) = w(\mathbf{s}, \mathbf{a}) + b$ .

For the operator  $\mathcal{T}^\xi$ , the two properties follows from

$$\begin{aligned} \mathcal{T}^\xi w(\mathbf{s}, \mathbf{a}) &= C(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_d[w(\mathbf{s}', \xi(\mathbf{s}'))] \geq C(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_d[w'(\mathbf{s}', \xi(\mathbf{s}'))] = \mathcal{T}^\xi w'(\mathbf{s}, \mathbf{a}), \quad \text{and} \\ \mathcal{T}^\xi(w + b)(\mathbf{s}, \mathbf{a}) &= C(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_d[(w(\mathbf{s}', \xi(\mathbf{s}')) + b)] = \mathcal{T}^\xi w(\mathbf{s}, \mathbf{a}) + \gamma b, \end{aligned}$$

It also follows from similar arguments that the operator  $\mathcal{T}^*$  also has the two properties, finishing the proof of assertion (i).

**Proof of Assertion (ii):** It follows from the contraction property of  $\mathcal{T}^\xi$  in Assertion (i) and the Banach Fixed Point Theorem (see Theorem 1.1 in [1]) that  $\mathcal{T}^\xi$  has a unique fixed point  $Q^\xi$ . To show that  $Q^\xi(\mathbf{s}, \xi(\mathbf{s})) = V_\alpha^\pi(\mathbf{s}; G)$  for induced policy  $\pi = \{\mathbf{a}(t)\}$  and all  $\mathbf{s}$ , temporarily define  $Q_V : \mathcal{SA}_\alpha \rightarrow \mathbb{R}$  by

$$Q_V(\mathbf{s}, \mathbf{a}) := C(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_d[V_\alpha^{\pi'}(f_\alpha(\mathbf{s}, \mathbf{a}, \mathbf{d}); G)],$$

for  $\pi' := \{\mathbf{a}(t+1) : t \geq 0\}$ . As a brief explanation,  $Q_V(\mathbf{s}, \mathbf{a})$  is the actor-value function if the system takes action  $\mathbf{a}$  in the first decision period and then follows policy  $\pi$ . The fact that  $\xi$  induces  $\pi = \{\mathbf{a}(t) : t \geq 0\}$  and the choice of  $\pi'$  imply  $\xi(\mathbf{s}) = \mathbf{a}(0)$  and

$$Q_V(\mathbf{s}, \xi(\mathbf{s})) := C(\mathbf{s}, \mathbf{a}(0)) + \gamma \mathbb{E}_d[V_\alpha^{\pi'}(f_\alpha(\mathbf{s}, \mathbf{a}(0), \mathbf{d}); G)] = V_\alpha^\pi(\mathbf{s}; G), \quad \text{for all } \mathbf{s} \in \mathcal{S}_\alpha. \quad (27)$$

Together with (24),

$$\mathcal{T}^\xi Q_V(\mathbf{s}, \mathbf{a}) = C(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_d[Q_V(\mathbf{s}', \xi(\mathbf{s}'))] = C(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_d[V_\alpha^\pi(\mathbf{s}'; G)] = Q_V(\mathbf{s}, \mathbf{a}),$$

for all  $(\mathbf{s}, \mathbf{a}) \in \mathcal{SA}_\alpha$  and  $\mathbf{s}' = f(\mathbf{s}, \mathbf{a}, \mathbf{d})$ . As a result,  $Q_V$  is a fixed point of  $\mathcal{T}^\pi$  and the uniqueness of  $Q^\xi$  implies  $Q_V = Q^\xi$ . Using (27) again,

$$Q^\xi(\mathbf{s}, \xi(\mathbf{s})) = Q_V(\mathbf{s}, \xi(\mathbf{s})) = V_\alpha^\pi(\mathbf{s}; G) \quad \text{for all } \mathbf{s} \in \mathcal{S}_\alpha.$$

**Proof of Assertion (iii):** It follows similar from the Banach fixed point theorem to the proof of Assertion (ii) that  $\mathcal{T}^*$  has a unique fixed point  $Q^*$ , where (25) immediately follows from Assertion (i). The fact that  $w_k \in \mathcal{C}$

for all  $\mathcal{C}$  follows immediately from Theorem 5. Because  $\mathcal{C}$  is a closed set with respect to norm  $\|\cdot\|_\infty$ , we obtain  $Q^* \in \mathcal{C}$  as the limit point of  $\{w_k\}$ .

**Proof of Assertion (iv):** It follows from Assertion (iii) that  $Q^*$  is convex, which is thus lower semi-continuous. Also, the set  $\mathcal{A}_\alpha(s)$  is closed for any  $s \in \mathcal{S}_\alpha$ . By Proposition 7.33 in [4] there exists a Markovian control  $\xi^* \in \Xi$  such that

$$\xi^*(s) \in \arg \min_{a \in \mathcal{A}_\alpha(s)} Q^*(s, a) \quad \text{for all } s \in \mathcal{S}_\alpha.$$

To show that  $Q^* = Q^{\xi^*}$ , for any  $(s, a) \in \mathcal{SA}_\alpha$  and  $s' = f(s, a, d)$ , we have

$$\begin{aligned} \mathcal{T}^{\xi^*}[Q^*(s, a)] &= C(s, a) + \gamma \mathbb{E}_d[Q^*(s', \xi^*(s'))] \\ &= C(s, a) + \gamma \mathbb{E}_d\left[\inf_{a' \in \mathcal{A}(s')} Q^*(s', a')\right] = \mathcal{T}^*Q^*(s, a) = Q^*(s, a), \end{aligned}$$

As a result,  $Q^*$  is a fixed point of  $\mathcal{T}^{\xi^*}$ , implying that  $Q^* = Q^{\xi^*}$ .

We next adopt existing results in [4], Proposition 9.8 to characterize  $V_\alpha^*(s; G)$  for all  $s \in \mathcal{S}_\alpha$ . First, the DSAA control problem is an infinite horizon stochastic optimal control problem as in Definition 9.1 in the reference under the discounted case. Apply Proposition 9.8, we obtain

$$V_\alpha^*(s; G) = \inf_{a \in \mathcal{A}_\alpha(s)} C(s, a) + \gamma E_{d \sim G}[V_\alpha^*(f_\alpha(s, a, d))], \quad \text{for all } s \in \mathcal{S}_\alpha.$$

Temporarily define

$$Q_V^*(s, a) := C(s, a) + \gamma E_{d \sim G}[V_\alpha^*(f_\alpha(s, a, d))], \quad \text{for all } (s, a) \in \mathcal{SA}_\alpha,$$

we can rewrite the equality by

$$\inf_{a \in \mathcal{A}_\alpha(s)} Q_V^*(s, a) = V_\alpha^*(s; G), \quad \text{for all } s \in \mathcal{S}_\alpha.$$

Together with (24), it holds for all  $(s, a) \in \mathcal{SA}_\alpha$  that

$$\mathcal{T}^*Q_V^*(s, a) = C(s, a) + \gamma E_{d \in \mathcal{G}}[V_\alpha^*(s'; G)] = Q_V^*(s, a), \quad \text{for } s' = f_\alpha(s, a, d).$$

Therefore,  $Q_V^* = Q^*$ , giving the second equality of (26). The first equality of (26) follows from

$$V_\alpha^{\pi^*}(s; G) = Q^{\xi^*}(s, \xi^*(s)) = Q^*(s, \xi^*(s)) = \min_{a \in \mathcal{A}_\alpha(s)} Q^*(s, a) = V_\alpha^*(s; G), \quad \text{for all } s \in \mathcal{S}.$$

Here, the first two equalities follow from Assertion (ii) and  $Q^* = Q^{\xi^*}$ , respectively.  $\square$

*Proof.* Proof of Lemma 6 and Corollary 2.

The statements in Lemma 6 immediately follow from Proposition 3(iii) and (iv). To prove Corollary 2, it follows from Proposition 3(iii) that  $Q^* \in \mathcal{C}$  is convex. By (iv),  $V_\alpha^*(\cdot; G)$  is convex.  $\square$

*Proof.* Proof of Corollary 1. Temporarily take  $L_0 = \bar{c}/(1 - \gamma)$ . We take  $w_0$  an  $L_0$ -ALC function on  $\mathcal{SA}_\alpha$  (for example, the zeroth function) and recursively generate the sequence  $\{w_k : k \geq 1\}$  as in Proposition 3(ii). By Theorem 4,  $w_k$  is  $L_0$ -ALC for all  $k \geq 1$ . Together with (25), we have shown that  $Q^*$  is  $L_0$ -ALC.

Next, Proposition 3(v) gives

$$V_\alpha^*(s; G) = \min_{a \in \mathcal{A}_\alpha(s)} Q^*(s, a), \quad \text{for all } s \in \mathcal{S}_\alpha.$$

For any  $s_1 \in \mathcal{S}$  and  $s_2 \in \mathcal{S}$ , Lemma 4 gives

$$\begin{aligned} V_\alpha^*(s_2; G) - V_\alpha^*(s_1; G) &= \min_{a_2 \in \mathcal{A}_\alpha(s_2)} \max_{a_1 \in \mathcal{A}_\alpha(s_1)} Q^*(s_2, a_2) - Q^*(s_1, a_1) \\ &\leq \min_{a_2 \in \mathcal{A}(s_2)} \max_{a_1 \in \mathcal{A}(s_1)} L_0 \rho_1(s_2 - s_1, a_2 - a_1) \leq L_0 \rho_0(s_2 - s_1). \end{aligned}$$

In particular, the first inequality follows from the definition of  $L_0$ -ALC and the second follows from Lemma 4. Finally, it is clear from the choice of  $\rho_0$  that, for  $s = (x, \bar{u}, R)$ ,

$$\rho_0(s) \leq (\bar{a} + 1) \left( \sum_{i=1}^n |x_i| + \sum_{j=1}^m |\bar{u}_j| + \sum_{\ell=1}^{L_n} \sum_{i=1}^n |R_{i\ell}| \right) \leq (\bar{a} + 1) p_S \|s\|_\infty.$$

As  $L_v := (\bar{a} + 1) p_S L_0 = \bar{c}(\bar{a} + 1) p_S / (1 - \gamma)$ , function  $V_\alpha^*(\cdot; G)$  is thus  $L_v$ -Lipschitz continuous.  $\square$

#### D.4 Proofs of Theorem 3.

In this section, we provide our proof of Theorem 3. Given Theorem 5, Theorem 4, and Proposition 3, most of the remaining arguments are standard, despite the proof is still long and technical. We hereby provide a roadmap and a list of supporting lemmas before we present the proof.

*Proof Roadmap.* The proof of Theorem 3 relies on an auxiliary Markovian policy and multiple versions of the value functions. Recall from Algorithm 1 that  $\xi_K$  is a measurable function from  $\mathcal{S}_\alpha$  to  $\mathbb{R}^{p_A}$ . The map  $\xi_K$  induces a Markovian policy  $\pi_K = \{\mathbf{a}(t)\}$ , for  $\mathbf{a}(t) = \xi_K(\mathbf{s}_\alpha(t))$  for all  $t \geq 0$  and  $\mathbf{s}_\alpha(t)$  from (23). Note that  $\pi_K$  is different from the output policy  $\hat{\pi}^*(\hat{D})$  of Algorithm 1 because policy  $\hat{\pi}^*(\hat{D})$  depends on  $\xi_K$  through the damped relation (3). For short, we write  $\hat{\pi}^*$  for  $\hat{\pi}^*(\hat{D})$ .

There are six different value functions with respect to three different transition dynamics. Take  $V^{\hat{\pi}^*}(\cdot; F)$  and  $V^*(\cdot; F)$  the value function of policy  $\hat{\pi}^*$  and the optimal value function under the true demand distribution, as defined by (8) and (9), respectively. Take  $V_{\alpha}^{\pi_K}(\cdot; F)$  and  $V_{\alpha}^*(\cdot; F)$  the value function of policy  $\pi_K$  and the optimal value function of damped control problem under the true demand distribution, respectively. Here, the controlled dynamics evolves by (23). Finally, take  $V_{\alpha}^{\pi_K}(\cdot; \hat{F})$  and  $V_{\alpha}^*(\cdot; \hat{F})$  the value function of policy  $\pi_K$  and the optimal value function of the DSAA control problem, respectively, where the demand follows from the sampling distribution and controlled dynamics evolve by (23). It holds that

$$\begin{aligned} 0 \leq V^{\hat{\pi}^*}(\mathbf{s}; F) - V^*(\mathbf{s}; F) &= (V^{\hat{\pi}^*}(\mathbf{s}; F) - V_{\alpha}^{\pi_K}(\mathbf{s}; F)) \\ &\quad + (V_{\alpha}^{\pi_K}(\mathbf{s}; F) - V_{\alpha}^{\pi_K}(\mathbf{s}; \hat{F})) + (V_{\alpha}^{\pi_K}(\mathbf{s}; \hat{F}) - V_{\alpha}^*(\mathbf{s}; \hat{F})) \\ &\quad + (V_{\alpha}^*(\mathbf{s}; \hat{F}) - V_{\alpha}^*(\mathbf{s}; F)) + (V_{\alpha}^*(\mathbf{s}; F) - V^*(\mathbf{s}; F)). \end{aligned} \quad (28)$$

We need separate estimates of the five terms on the right-hand side of (28). Intuitively, the first and last terms on the right-hand side of (28) capture the error of using a damped system to approximate the nond-damped system. To bound the two terms, we employ coupling arguments in Section D.6 that align the damped system and a non-damped system in the same probability space and bound the difference of the cost function on each sample path. The second and fourth terms on the right-hand side of (28) capture the error of using a sampling distribution to approximate the unknown true distribution. To bound this source of error, we employ covering arguments in Section D.7 on the compact domain  $\mathcal{S}_\alpha$  and then make use of an existing concentration result. Finally, the third term on the right-hand side of (28) gives the error of having a finite number of iterations in the CFQI algorithm. This source of error can be bounded by the contraction of the Bellman operator  $\mathcal{T}^*$ , see Section D.5. The following four supporting lemmas provide an upper bound for each of the five terms.

**Lemma 6.** For any  $\alpha \in (0, 1)$  and  $\mathbf{s} \in \mathcal{S}_\alpha$ ,

$$0 \leq \sup_{\mathbf{s} \in \mathcal{S}_\alpha} (V_{\alpha}^{\pi_K}(\mathbf{s}; \hat{F}) - V_{\alpha}^*(\mathbf{s}; \hat{F})) = \frac{4\varepsilon_m}{(1-\gamma)^2}. \quad (29)$$

**Lemma 7.** Take  $\alpha \in (0, 1)$  and

$$\mathcal{E}_1(\delta) := \frac{3(\bar{a}+1)\bar{c}p_S}{(1-\gamma)^2} \sqrt{\frac{p_S + p_A}{2M} \log \frac{3M(\bar{a}+1)x_\alpha}{\delta \bar{d}}}.$$

For any  $\pi \in \Pi_\alpha$ , event

$$\|V_{\alpha}^{\pi}(\cdot; F) - V_{\alpha}^{\pi}(\cdot; \hat{F})\|_{\infty} \leq \frac{\gamma \|V_{\alpha}^{\pi}(\cdot; \hat{F}) - V_{\alpha}^*(\cdot; \hat{F})\|_{\infty}}{1-\gamma} + \mathcal{E}_1(\delta) \quad (30)$$

has probability at least  $1 - \delta$ . Also, event

$$\|V_{\alpha}^{\pi}(\cdot; F) - V_{\alpha}^{\pi}(\cdot; \hat{F})\|_{\infty} \leq \frac{\gamma \|V_{\alpha}^{\pi}(\cdot; F) - V_{\alpha}^*(\cdot; F)\|_{\infty}}{1-\gamma} + \mathcal{E}_1(\delta) \quad (31)$$

has probability at least  $1 - \delta$ . Here, the probability is taken with respect to the sample  $\hat{D}$ .

**Lemma 8.** For any  $\alpha \in (0, 1)$  and  $\mathbf{s} \in \mathcal{S}_\alpha$ ,

$$V^{\hat{\pi}^*}(\mathbf{s}; F) - V_{\alpha}^{\pi_K}(\mathbf{s}; F) \leq \frac{\alpha}{(1-\gamma)} V_{\alpha}^{\pi_K}(\mathbf{s}; F). \quad (32)$$

**Lemma 9.** For any  $\alpha \in (0, 1)$  and  $\mathbf{s} \in \mathcal{S}_\alpha$ ,

$$V_{\alpha}^*(\mathbf{s}; F) - V^*(\mathbf{s}; F) \leq \frac{\alpha \bar{c} m (L_n + 1) \bar{d}}{(1-\gamma)^2}. \quad (33)$$

*Proof.* Proof of Theorem 3. First, the choice of  $\alpha$  in (5) indicates that  $\mathbf{s} \in \mathcal{S}_\alpha$ , so that Lemmas 6 – 9 are all applicable. Similar to the treatment in (28), we can expand the right-hand side of Lemma 8 by

$$\begin{aligned} \frac{\alpha}{1-\gamma} V_\alpha^{\pi_K}(\mathbf{s}; F) &= \frac{\alpha}{1-\gamma} V^*(\mathbf{s}; F) + \frac{\alpha}{1-\gamma} (V_\alpha^{\pi_K}(\mathbf{s}; F) - V_\alpha^{\pi_K}(\mathbf{s}; \hat{F})) + \frac{\alpha}{1-\gamma} (V_\alpha^{\pi_K}(\mathbf{s}; \hat{F}) - V_\alpha^*(\mathbf{s}; \hat{F})) \\ &\quad + \frac{\alpha}{1-\gamma} (V_\alpha^*(\mathbf{s}; \hat{F}) - V_\alpha^*(\mathbf{s}; F)) + \frac{\alpha}{1-\gamma} (V_\alpha^*(\mathbf{s}; F) - V^*(\mathbf{s}; F)). \end{aligned}$$

Combine (28) with (32) and using  $\alpha \leq (1-\gamma)/\sqrt{M} \leq 1-\gamma$ , we obtain

$$\begin{aligned} 0 &\leq V^{\hat{\pi}^*}(\mathbf{s}; F) - V^*(\mathbf{s}; F) \\ &\leq \frac{1}{\sqrt{M}} V^*(\mathbf{s}; F) + 2|V_\alpha^{\pi_K}(\mathbf{s}; F) - V_\alpha^{\pi_K}(\mathbf{s}; \hat{F})| + 2|V_\alpha^{\pi_K}(\mathbf{s}; \hat{F}) - V_\alpha^*(\mathbf{s}; \hat{F})| \\ &\quad + 2(V_\alpha^*(\mathbf{s}; \hat{F}) - V_\alpha^*(\mathbf{s}; F)) + 2(V_\alpha^*(\mathbf{s}; F) - V^*(\mathbf{s}; F)), \end{aligned} \quad (34)$$

In Lemma 7, we take  $\pi = \pi_K \in \Pi_\alpha$ ,

$$|V_\alpha^{\pi_K}(\mathbf{s}; F) - V_\alpha^{\pi_K}(\mathbf{s}; \hat{F})| \leq \frac{\gamma \|V_\alpha^{\pi_K}(\cdot; \hat{F}) - V_\alpha^*(\cdot; \hat{F})\|_\infty}{1-\gamma} + \mathcal{E}_1(\delta/2)$$

holds with probability at least  $1 - \delta/2$ . Together with (29) in Lemma 6,

$$|V_\alpha^{\pi_K}(\mathbf{s}; F) - V_\alpha^{\pi_K}(\mathbf{s}; \hat{F})| \leq \frac{4\gamma\varepsilon_m}{(1-\gamma)^3} + \mathcal{E}_1(\delta/2) \quad (35)$$

holds with probability at least  $1 - \delta/2$ .

In Lemma 7, we next take  $\pi$  such that  $V_\alpha^\pi(\mathbf{s}; F) = V_\alpha^*(\mathbf{s}; F)$ , the existence of such  $\pi$  follows from Proposition 3(v). We have

$$V_\alpha^*(\mathbf{s}; \hat{F}) - V_\alpha^*(\mathbf{s}; F) \leq V_\alpha^\pi(\mathbf{s}; \hat{F}) - V_\alpha^\pi(\mathbf{s}; F).$$

Together with (31),

$$V_\alpha^*(\mathbf{s}; \hat{F}) - V_\alpha^*(\mathbf{s}; F) \leq \mathcal{E}_1(\delta/2) \quad (36)$$

holds with probability at least  $1 - \delta/2$ .

Plug (35), (29), (36), and (33) in Lemma 9 to (34), we can compute that

$$\begin{aligned} 0 &\leq V^{\hat{\pi}^*}(\mathbf{s}; F) - V^*(\mathbf{s}; F) \\ &\leq \frac{1}{\sqrt{M}} V^*(\mathbf{s}; F) + \left(2\mathcal{E}_1(\delta/2) + \frac{8\gamma\varepsilon_m}{(1-\gamma)^3}\right) + \frac{8\varepsilon_m}{(1-\gamma)^2} + 2\mathcal{E}_1(\delta/2) + \frac{2\alpha\bar{c}m(L_n+1)\bar{d}}{(1-\gamma)^2} \\ &\leq \frac{1}{\sqrt{M}} V^*(\mathbf{s}; F) + 4\mathcal{E}_1(\delta/2) + \frac{8\varepsilon_m}{(1-\gamma)^3} + \frac{2\alpha\bar{c}m(L_n+1)\bar{d}}{(1-\gamma)^2} \end{aligned}$$

with probability at least  $1 - \delta$ , where the last inequality utilizes the choice that  $\alpha \leq M^{-1}(1-\gamma)$ . We take

$$\mathcal{E}(\delta) := 4\mathcal{E}_1(\delta/2) + \frac{2\alpha\bar{c}m(L_n+1)\bar{d}}{(1-\gamma)^2},$$

giving (5). Finally, using the expression of  $\mathcal{E}_1(\delta)$  and  $x_\alpha$ , it is clear to check that

$$\mathcal{E}(\delta) = O\left(\frac{\bar{a}\bar{c}(p_S + p_A)^2}{(1-\gamma)^2} \sqrt{\frac{1}{M} \log \frac{\bar{a}M}{\alpha\delta}}\right), \quad \text{for } \alpha = \min\left\{\frac{1-\gamma}{\sqrt{M}}, \frac{\bar{d}}{\|\mathbf{s}\|_\infty}\right\},$$

where  $O(\cdot)$  hides an “absolute” constant that does not depend on any model primitives.  $\square$

## D.5 The Contraction Arguments

The proof of Lemma 6 mainly utilizes the contraction property of the Bellman operator  $\mathcal{T}^*$ . The main line of the proof is similar to [13], while the presence of the convex optimization penalty  $\Lambda$  gives an additional source of difficulty.

*Proof of Lemma 6.* In this proof, we fix the damping factor  $\alpha$  and  $G = \hat{F}$  in Proposition 3. Note that the inequality in (29) is trivial and our goal is to prove the equality. Recall  $\xi_K$  and  $Q_K$  from Algorithm 1. It is clear that

$$Q_K(\mathbf{s}, \xi_K(\mathbf{s})) \leq Q_K(\mathbf{s}, \xi(\mathbf{s})) \quad \text{and} \quad \mathcal{T}^{\xi_K} Q_K \leq \mathcal{T}^\xi Q_K \quad \text{for all } \mathbf{s} \in \mathcal{S} \quad \text{and} \quad \xi \in \Xi. \quad (37)$$

Because  $\pi_K$  is the induced policy by  $\xi_K$ , it follows from Proposition 3(ii) that

$$V_\alpha^{\pi_K}(\mathbf{s}; \hat{F}) = Q^{\xi_K}(\mathbf{s}, \xi_K(\mathbf{s})) \quad \text{for all } \mathbf{s} \in \mathcal{S}_\alpha.$$

Next, recall from Proposition 3(v) the optimal Markovian policy  $\xi^*$ . We similarly obtain from Proposition 3(ii) that

$$V_\alpha^*(\mathbf{s}; \widehat{F}) = Q^{\xi^*}(\mathbf{s}, \xi^*(\mathbf{s})) \quad \text{for all } \mathbf{s} \in \mathcal{S}_\alpha.$$

We can now compute that

$$\begin{aligned} & V_\alpha^{\pi_K}(\mathbf{s}; \widehat{F}) - V_\alpha^*(\mathbf{s}; \widehat{F}) \\ &= (Q^{\xi_K}(\mathbf{s}, \xi_K(\mathbf{s})) - Q_K(\mathbf{s}, \xi_K(\mathbf{s}))) + (Q_K(\mathbf{s}, \xi_K(\mathbf{s})) - Q_K(\mathbf{s}, \xi^*(\mathbf{s}))) \\ & \quad + (Q_K(\mathbf{s}, \xi^*(\mathbf{s})) - Q^{\xi^*}(\mathbf{s}, \xi^*(\mathbf{s}))) \\ &\leq (Q^{\xi_K}(\mathbf{s}, \xi_K(\mathbf{s})) - Q_K(\mathbf{s}, \xi_K(\mathbf{s}))) + (Q_K(\mathbf{s}, \xi^*(\mathbf{s})) - Q^{\xi^*}(\mathbf{s}, \xi^*(\mathbf{s}))) \\ &\leq \|Q^{\xi_K} - Q_K\|_\infty + \|Q_K - Q^{\xi^*}\|_\infty \\ &\leq \|Q^{\xi_K} - Q^{\xi^*}\|_\infty + 2\|Q_K - Q^{\xi^*}\|_\infty, \quad \text{for all } \mathbf{s} \in \mathcal{S}. \end{aligned} \tag{38}$$

Here, the first inequality follows from taking  $\xi = \xi^*$  in the first inequality of (37). The second and third inequalities are obtained by taking supreme over  $\mathcal{S}\mathcal{A}_\alpha$  and then adopting the triangular inequality. Next, it follows from Proposition 3(ii) and (iii) that  $Q^{\xi_K} = \mathcal{T}^{\xi_K} Q^{\xi_K}$  and  $Q^{\xi^*} = \mathcal{T}^{\xi^*} Q^{\xi^*}$ . We can compute that

$$\begin{aligned} 0 &\leq Q^{\xi_K} - Q^{\xi^*} \\ &= (\mathcal{T}^{\xi_K} Q^{\xi_K} - \mathcal{T}^{\xi_K} Q^{\xi^*}) + (\mathcal{T}^{\xi_K} Q^{\xi^*} - \mathcal{T}^{\xi_K} Q_K) + (\mathcal{T}^{\xi_K} Q_K - \mathcal{T}^{\xi^*} Q_K) + (\mathcal{T}^{\xi^*} Q_K - \mathcal{T}^{\xi^*} Q^{\xi^*}) \\ &\leq \|\mathcal{T}^{\xi_K} Q^{\xi_K} - \mathcal{T}^{\xi_K} Q^{\xi^*}\|_\infty + \|\mathcal{T}^{\xi_K} Q^{\xi^*} - \mathcal{T}^{\xi_K} Q_K\|_\infty + \|\mathcal{T}^{\xi^*} Q_K - \mathcal{T}^{\xi^*} Q^{\xi^*}\|_\infty \\ &\leq \gamma\|Q^{\xi_K} - Q^{\xi^*}\|_\infty + 2\gamma\|Q^{\xi^*} - Q_K\|_\infty, \end{aligned}$$

where the first inequality follows from taking  $\xi = \xi^*$  in (37) and the second inequality follows from Proposition 3(i). As a result,

$$\|Q^{\xi_K} - Q^{\xi^*}\|_\infty \leq \frac{2\gamma}{1-\gamma}\|Q^{\xi^*} - Q_K\|_\infty.$$

Together with (38),

$$\sup_{\mathbf{s} \in \mathcal{S}} (V_\alpha^{\pi_K}(\mathbf{s}; \widehat{F}) - V_\alpha^*(\mathbf{s}; \widehat{F})) \leq \frac{2}{1-\gamma}\|Q^{\xi^*} - Q_K\|_\infty. \tag{39}$$

To estimate the right-hand side of (39), recall  $Q_k$  from Algorithm 1. For  $k = [K] \cup \{0\}$ ,  $\mathbf{s} \in \mathcal{S}_\alpha$  and  $\mathbf{a} \in \mathcal{A}_\alpha(\mathbf{s})$ , let

$$Q_k^\Lambda(\mathbf{s}, \mathbf{a}) := Q_k(\mathbf{s}, \mathbf{a}) + \Lambda(\mathbf{s}, \mathbf{a}; \eta^k) \quad \text{and} \quad \xi_k(\mathbf{s}) \in \arg \min_{\mathbf{a} \in \mathcal{A}(\mathbf{s})} Q_k^\Lambda(\mathbf{s}, \mathbf{a}). \tag{40}$$

By (13) and (24), it holds that

$$\mathcal{T}^{\xi^*} Q_k^\Lambda \geq \mathcal{T}^* Q_k^\Lambda = \mathcal{T}^{\xi_k} Q_k^\Lambda, \quad \text{for all } k \in [K] \cup \{0\}. \tag{41}$$

Together with  $Q^{\xi^*} = \mathcal{T}^{\xi^*} Q^{\xi^*}$ , we can compute that

$$\begin{aligned} Q^{\xi^*} - Q_{k+1}^\Lambda &= (\mathcal{T}^{\xi^*} Q^{\xi^*} - \mathcal{T}^{\xi^*} Q_k^\Lambda) + (\mathcal{T}^{\xi^*} Q_k^\Lambda - \mathcal{T}^{\xi_k} Q_k^\Lambda) + (\mathcal{T}^{\xi_k} Q_k^\Lambda - Q_{k+1}^\Lambda) \\ &\geq \mathcal{T}^{\xi^*} Q^{\xi^*} - \mathcal{T}^{\xi^*} Q_k^\Lambda + \xi_k, \end{aligned}$$

where we take  $\epsilon_k := \mathcal{T}^{\xi_k} Q_k^\Lambda - Q_{k+1}^\Lambda$ , for  $k \in [K-1] \cup \{0\}$ . Together with Proposition 3(i),

$$Q^{\xi^*} - Q_{k+1}^\Lambda \geq -\gamma\|Q^{\xi^*} - Q_k^\Lambda\|_\infty - \|\epsilon_k\|_\infty.$$

It similarly follows from (24) that  $\mathcal{T}^{\xi^*} Q^{\xi^*} \leq \mathcal{T}^{\xi_k} Q^{\xi^*}$  for all  $k \in [K] \cup \{0\}$ , implying that

$$Q^{\xi^*} - Q_{k+1}^\Lambda = (\mathcal{T}^{\xi^*} Q^{\xi^*} - \mathcal{T}^{\xi_k} Q^{\xi^*}) + (\mathcal{T}^{\xi_k} Q^{\xi^*} - \mathcal{T}^{\xi_k} Q_k^\Lambda) + \epsilon_k \leq \mathcal{T}^{\xi_k} Q^{\xi^*} - \mathcal{T}^{\xi_k} Q_k^\Lambda + \epsilon_k,$$

for  $k \leq K-1$ . Together with Proposition 3(i),

$$Q^{\xi^*} - Q_{k+1}^\Lambda \leq \gamma\|Q^{\xi^*} - Q_k^\Lambda\|_\infty + \|\epsilon_k\|_\infty.$$

In particular, we have shown that

$$\|Q^{\xi^*} - Q_{k+1}^\Lambda\|_\infty \leq \gamma\|Q^{\xi^*} - Q_k^\Lambda\|_\infty + \|\epsilon_k\|_\infty, \quad \text{for all } k \in [K-1] \cup \{0\}.$$

By recursion,

$$\|Q^{\xi^*} - Q_K^\Lambda\|_\infty \leq \sum_{k=0}^{K-1} \gamma^{K-k-1} \|\epsilon_k\|_\infty + \gamma^K \|Q^{\xi^*} - Q_0^\Lambda\|_\infty.$$

By Corollary 2,  $Q^{\xi^*} \in \mathcal{C}$ . Together with  $Q_0 \in \mathcal{C}$  and  $\|\Lambda(\cdot; \eta^k)\|_\infty \leq \eta^k B_\alpha$  for all  $k \geq 0$ , it holds that

$$-2B_\alpha \leq -Q_0 - B_\alpha \leq Q^{\xi^*} - Q_0^\Lambda \leq Q^{\xi^*} + B_\alpha \leq 2B_\alpha.$$

Therefore,  $\|Q^{\xi^*} - Q_0^\Lambda\|_\infty \leq 2B_\alpha$ . Together with  $\|Q_K^\Lambda - Q_K\|_\infty \leq \eta^K B_\alpha$ ,

$$\|Q^{\xi^*} - Q_K\|_\infty \leq \sum_{k=0}^{K-1} \gamma^{K-k-1} \|\epsilon_k\|_\infty + 2B_\alpha \gamma^K + B_\alpha \eta^K.$$

Finally, the choice of  $Q^k$  as in Algorithm 1, line 2 gives  $\|Q_{k+1} - \mathcal{T}^* Q_k^\Lambda\|_\infty \leq \varepsilon_m$ . Together with (41) and (40), we obtain

$$\|\epsilon_k\|_\infty \leq \|Q_{k+1} - \mathcal{T}^{\xi_k} Q_k^\Lambda\|_\infty + \|\Lambda(\cdot; \eta^{k+1})\|_\infty \leq \varepsilon_m + \eta^k B_\alpha,$$

for all  $j \in [K]$ . Therefore, for  $\eta < \gamma < 1$ , we obtain

$$\sum_{k=0}^{K-1} \gamma^{K-k-1} \|\epsilon_k\|_\infty \leq \sum_{k=0}^{K-1} \gamma^{K-k-1} (\varepsilon_m + \eta^k B_\alpha) = \frac{\varepsilon_m(1 - \gamma^K)}{1 - \gamma} + \frac{(\gamma^K - \eta^K)B_\alpha}{\gamma - \eta}.$$

It follows from  $\eta = \max\{\gamma/2, 2\gamma - 1\} < \gamma$  that  $\gamma - \eta = \min\{\gamma/2, 1 - \gamma\} \leq \gamma(1 - \gamma)/2$ , and

$$2\gamma^K + \eta^K + \frac{\gamma^K - \eta^K}{\gamma - \eta} \leq 2\gamma^K + \frac{\gamma^K}{\gamma - \eta} \leq \frac{2\gamma^{K-1}}{1 - \gamma} + \frac{\gamma^K}{\gamma/2(1 - \gamma)} \leq \frac{4\gamma^{K-1}}{1 - \gamma},$$

implying that

$$\|Q^{\xi^*} - Q_K\|_\infty \leq (\varepsilon_m + 4\gamma^{K-1}B_\alpha)/(1 - \gamma).$$

Recall that

$$B_\alpha = \frac{\bar{c}(m+n)x_\alpha}{1 - \gamma} = \frac{3\bar{a}\bar{c}\bar{d}L_n(m+n)}{\alpha^2(1 - \gamma)}.$$

The choice of  $K$  gives  $4\gamma^{K-1}B_\alpha \leq \varepsilon_m$ . Together with (39), we have shown (29), finishing the proof.

## D.6 The Coupling Arguments

In this section, we provide coupling arguments to prove Lemmas 8 and 9. The proof of Proposition 2 also appears in this section. Throughout the section, all demands are sampled from the unknown distribution  $F$ .

*Proof of Lemma 8.* We consider the following coupling between two systems, both with initial state  $\mathbf{s}(0)$ . System 1 is controlled by policy  $\hat{\pi}^* = \{\hat{\mathbf{a}}(t) : t \geq 0\}$  and has no damping, where the controlled dynamics  $\{\hat{\mathbf{s}}(t)\}$  evolves as

$$\hat{\mathbf{s}}(t+1) = f(\hat{\mathbf{s}}(t), \hat{\mathbf{a}}(t), \mathbf{d}(t+1)) \quad \text{for } \hat{\mathbf{s}}(t) = (\hat{\mathbf{x}}(t), \hat{\mathbf{u}}(t), \hat{R}(t)), \quad t \geq 0.$$

System 2 has damping factor  $\alpha$  and is controlled by policy  $\pi_K = \{\mathbf{a}(t) : t \geq 0\}$ , where  $\mathbf{a}(t) = \xi_K(\mathbf{s}(t))$  and the controlled state  $\{\mathbf{s}(t)\}$  evolves as

$$\mathbf{s}(t+1) = f_\alpha(\mathbf{s}(t), \mathbf{a}(t), \mathbf{d}(t+1)), \quad \text{for } \mathbf{s}(t) = (\mathbf{x}(t), \tilde{\mathbf{u}}(t), R(t)), \quad t \geq 0.$$

We can compute that

$$V^{\hat{\pi}^*}(\mathbf{s}; F) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t C(\hat{\mathbf{s}}(t), \hat{\mathbf{a}}(t))\right] \quad \text{and} \quad V_\alpha^{\pi_K}(\mathbf{s}; F) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t C(\mathbf{s}(t), \mathbf{a}(t))\right],$$

and thus

$$V^{\hat{\pi}^*}(\mathbf{s}; F) - V_\alpha^{\pi_K}(\mathbf{s}; F) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t C(\hat{\mathbf{s}}(t) - \mathbf{s}(t), \hat{\mathbf{a}}(t) - \mathbf{a}(t))\right]. \quad (42)$$

We next show that

$$\begin{aligned} \hat{\mathbf{s}}(t) - \mathbf{s}(t) &= \alpha \sum_{\ell=0}^{t-1} (\mathbf{x}(\ell) - A^\top \mathbf{y}(\ell), \tilde{\mathbf{u}}(\ell) - \mathbf{y}(\ell), O) \\ \mathbf{r}(t) &= \hat{\mathbf{s}}(t) - \mathbf{s}(t), \quad \text{and} \quad \hat{\mathbf{a}}(t) = \mathbf{a}(t) \quad \text{for all } t \geq 0, \end{aligned} \quad (43)$$

where the sequence  $\mathbf{r}(t)$  is defined in (3). Proof by induction. When  $t = 0$ , it is clear that  $\mathbf{s}(0) = \hat{\mathbf{s}}(0)$  and (3) implies that  $\hat{\mathbf{a}}(0) = \mathbf{a}(0)$ . Now take the induction hypothesis that (43) holds for some  $t_0$  and  $t = t_0 + 1$ . For

the first equality in (43), we can compute that

$$\begin{aligned}
\widehat{\mathbf{s}}(t) &= f(\widehat{\mathbf{s}}(t_0), \widehat{\mathbf{a}}(t_0), \mathbf{d}(t_0 + 1)) \\
&= f\left(\mathbf{s}(t_0) + \alpha \sum_{\ell=0}^{t_0-1} (\mathbf{x}(\ell) - A^\top \mathbf{y}(\ell), \tilde{\mathbf{u}}(\ell) - \mathbf{y}(\ell), O), \mathbf{a}(t_0), \mathbf{d}(t_0 + 1)\right) \\
&= f(\mathbf{s}(t_0), \mathbf{a}(t_0), \mathbf{d}(t_0 + 1)) + \alpha \sum_{\ell=0}^{t_0-1} f\left((\mathbf{x}(\ell) - A^\top \mathbf{y}(\ell), \tilde{\mathbf{u}}(\ell) - \mathbf{y}(\ell), O), \mathbf{0}, \mathbf{0}\right) \\
&= f_\alpha(\mathbf{s}(t_0), \mathbf{a}(t_0), \mathbf{d}(t_0 + 1)) + \alpha \sum_{\ell=0}^{t_0-1} (\mathbf{x}(\ell) - A^\top \mathbf{y}(\ell), \tilde{\mathbf{u}}(\ell) - \mathbf{y}(\ell), O) \\
&= \mathbf{s}(t_0 + 1) + \alpha \sum_{\ell=0}^{t_0-1} (\mathbf{x}(\ell) - A^\top \mathbf{y}(\ell), \tilde{\mathbf{u}}(\ell) - \mathbf{y}(\ell), O) \\
&= \mathbf{s}(t) + \alpha \sum_{\ell=0}^{t-1} (\mathbf{x}(\ell) - A^\top \mathbf{y}(\ell), \tilde{\mathbf{u}}(\ell) - \mathbf{y}(\ell), O).
\end{aligned}$$

Here, note that the second equality is the induction hypothesis, the third equality uses the linearity of  $f$ , the fourth equality uses the fact that

$$f(\mathbf{s}, \mathbf{a}, \mathbf{d}) = f_\alpha(\mathbf{s}, \mathbf{a}, \mathbf{d}) + \alpha(\mathbf{x} - A^\top \mathbf{y}, \tilde{\mathbf{u}} - \mathbf{y}, O) \quad \text{and} \quad f(\mathbf{s}, \mathbf{0}, \mathbf{0}) = \mathbf{s},$$

for all  $\mathbf{s} = (\mathbf{x}, \tilde{\mathbf{u}}, \mathbf{d})$  and  $\mathbf{a} = (\mathbf{y}, \mathbf{z})$ . As a result, the first equality in (43) holds for  $t = t_0 + 1$ , which also implies that

$$\alpha(\mathbf{x}(t_0) - A^\top \mathbf{y}(t_0), \tilde{\mathbf{u}}(t_0) - \mathbf{y}(t_0), O) = \widehat{\mathbf{s}}(t) - \mathbf{s}(t) - \widehat{\mathbf{s}}(t_0) + \mathbf{s}(t_0). \quad (44)$$

To show that the second equality in (43) holds for  $t = t_0 + 1$ , observe that the induction hypothesis gives

$$\widehat{\mathbf{y}}(t_0) = \mathbf{y}(t_0) \quad \text{and} \quad \mathbf{r}(t_0) = \widehat{\mathbf{s}}(t_0) - \mathbf{s}(t_0).$$

Using the choice of  $\mathbf{r}(t)$  in (3),

$$\begin{aligned}
\mathbf{r}(t) &= (1 - \alpha)\mathbf{r}(t_0) + \alpha(\widehat{\mathbf{x}}(t_0) - A\widehat{\mathbf{y}}(t_0), \widehat{\mathbf{u}}(t_0) - \widehat{\mathbf{y}}(t_0), O) \\
&= (1 - \alpha)\widehat{\mathbf{s}}(t_0) - \mathbf{s}(t_0) + \alpha(\mathbf{x}(t_0) - A\mathbf{y}(t_0), \tilde{\mathbf{u}}(t_0) - \mathbf{y}(t_0), O) + \alpha(\widehat{\mathbf{s}}(t_0) - \mathbf{s}(t_0)) \\
&= \widehat{\mathbf{s}}(t) - \mathbf{s}(t).
\end{aligned}$$

Here, the first equality follows from the induction hypothesis and the second equality utilizes (44). As a result, the second equality in (43) holds for  $t = t_0 + 1$ . Also, the choice of  $\widehat{\mathbf{a}}(t)$  in (3) gives

$$\widehat{\mathbf{a}}(t) = \xi_K(\widehat{\mathbf{s}}(t) - \mathbf{r}(t)) = \xi_K(\mathbf{s}(t)) = \mathbf{a}(t), \quad \text{for } t = t_0 + 1.$$

To summarize, we have shown that (43) holds for  $t = t_0 + 1$ . By induction, (43) holds for all  $t \geq 0$ .

We can now compute that, for all  $t \geq 0$ ,

$$\begin{aligned}
C(\widehat{\mathbf{s}}(t) - \mathbf{s}(t), \widehat{\mathbf{a}}(t) - \mathbf{a}(t)) &= C\left(\alpha \sum_{\ell=0}^{t-1} (\mathbf{x}(\ell) - A^\top \mathbf{y}(\ell), \tilde{\mathbf{u}}(\ell) - \mathbf{y}(\ell), O), \mathbf{0}\right) \\
&= \alpha \sum_{\ell=0}^{t-1} \mathbf{h}^\top (\mathbf{x}(\ell) - A^\top \mathbf{y}(\ell)) + \mathbf{b}^\top (\tilde{\mathbf{u}}(\ell) - \mathbf{y}(\ell)) \\
&= \alpha \sum_{\ell=0}^{t-1} C(\mathbf{s}(\ell), \mathbf{a}(\ell)).
\end{aligned}$$

Notice that  $C(\mathbf{s}(\ell), \mathbf{a}(\ell)) \geq 0$  holds for all  $t$  due to the expression of  $C$  in (7) and  $\mathbf{a}(t) \in \mathcal{A}_\alpha(\mathbf{s}(t))$ . Together with (42), we have shown that

$$\begin{aligned}
V^{\widehat{\pi}^*}(\mathbf{s}; F) - V_{\alpha}^{\pi_K}(\mathbf{s}; F) &= \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t C(\widehat{\mathbf{s}}(t) - \mathbf{s}(t), \widehat{\mathbf{a}}(t) - \mathbf{a}(t))\right] \\
&= \alpha \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t \sum_{\ell=0}^{t-1} C(\mathbf{s}(\ell), \mathbf{a}(\ell))\right] \\
&= \alpha \mathbb{E}\left[\sum_{\ell=0}^{\infty} C(\mathbf{s}(\ell), \mathbf{a}(\ell)) \left(\sum_{t=\ell+1}^{\infty} \gamma^t\right)\right] \\
&\leq \frac{\alpha}{1-\gamma} \mathbb{E}\left[\sum_{\ell=0}^{\infty} C(\mathbf{s}(\ell), \mathbf{a}(\ell)) \gamma^\ell\right] = \frac{\alpha}{1-\gamma} V_{\alpha}^{\pi_K}(\mathbf{s}; F), \quad (45)
\end{aligned}$$

finishing the proof.  $\square$

*Proof of Proposition 2.* Assertion (i) follows from Proposition 7.33 in [3] and that  $Q_K$  is a convex function.

We next prove Assertion (ii). Using the induction in the proof of Lemma 8, we have shown that  $\widehat{\mathbf{s}}(t) - \mathbf{r}(t) = \mathbf{s}(t)$  and  $\widehat{\mathbf{a}}(t) = \mathbf{a}(t)$  for all  $t$ . It is known that  $\{\mathbf{a}(t)\}$  is a feasible policy for the damped control problem, that is,  $\mathbf{a}(t) \in \mathcal{A}_\alpha(\mathbf{s}(t))$  and  $\mathbf{a}(t) \in \mathcal{F}_t$ . As a result,

$$\widehat{\mathbf{a}}(t) \in \mathcal{A}_\alpha(\widehat{\mathbf{s}}(t) - \mathbf{r}(t)).$$

By recursion, it is easy to show that, for any  $t \geq 0$ ,  $\widehat{\mathbf{s}}(t)$ ,  $\mathbf{r}(t)$ ,  $\widehat{\mathbf{a}}(t) = \mathbf{a}(t)$  are all independent of  $(\mathbf{d}(t+1), \mathbf{d}(t+2), \dots)$ . As a result,  $\{\widehat{\mathbf{a}}(t)\}$  is an admissible policy.  $\square$

*Proof of Lemma 9.* First, recall that

$$V^*(\mathbf{s}; F) = \inf_{\pi \in \Pi} V^\pi(\mathbf{s}; F).$$

We consider a subfamily of  $\Pi_S \subseteq \Pi$  such that any policy  $\pi = \{\mathbf{a}(t)\} \in \Pi_S$ , in period  $t$ , can order a maximum number of components that is just sufficient to assemble products that meet the unsatisfied demand at period  $t$ , plus  $L_n$  additional periods of maximum demand. Formally, we require  $\mathbf{z}(t) \leq A^\top(\tilde{\mathbf{u}}(t) + L_n \bar{\mathbf{d}}\mathbf{1})$ , for all  $t \geq 0$ . As there is no need at period  $t$  to order components for assembly more than  $t + L_n$  periods later, and thus there is no loss of optimality by restricting to the subfamily  $\Pi_S$ , namely,

$$V^*(\mathbf{s}; F) = \inf_{\pi \in \Pi_S} V^\pi(\mathbf{s}; F).$$

Our goal is to show that, for any policy  $\pi \in \Pi_S$ , there is a corresponding policy  $\pi_\alpha \in \Pi_\alpha$ , such that

$$V_\alpha^{\pi_\alpha}(\mathbf{s}; F) \leq V^\pi(\mathbf{s}; F) + \frac{\alpha \bar{c} m (L_n + 1) \bar{d}}{(1 - \gamma)^2}. \quad (46)$$

Once (46) is proved, the lemma follows from taking the infimum over  $\pi \in \Pi_S$ .

We again consider a coupling between two systems, both with initial state  $\mathbf{s}(0)$ . System 1 is not damped and controlled by a fixed policy  $\pi = \{\mathbf{a}(t)\} \in \Pi_S$ . Recall from (3) that  $\{\mathbf{s}(t)\}$  is the controlled state of system 1, which is specified by (6). System 2 has damping factor  $\alpha$  and is controlled by policy  $\pi_\alpha = \{\mathbf{a}_\alpha(t) : t \geq 0\}$  that we specify later. Denote by  $\{\mathbf{s}_\alpha(t)\}$  the controlled state of system 2, namely,

$$\mathbf{s}_\alpha(t+1) = f_\alpha(\mathbf{s}_\alpha(t), \mathbf{a}_\alpha(t), \mathbf{d}(t+1)), \quad \text{for } \mathbf{s}_\alpha(t) = (\mathbf{x}_\alpha(t), \tilde{\mathbf{u}}_\alpha(t), R_\alpha(t)), \quad \text{for all } t \geq 0.$$

To specify the coupling between  $\mathbf{a}(t)$  and  $\mathbf{a}_\alpha(t)$ , take

$$\mathbf{a}(t) = (\mathbf{y}(t), \mathbf{z}(t)) \quad \text{and} \quad \mathbf{a}_\alpha(t) = \{\mathbf{y}_\alpha(t), \mathbf{z}_\alpha(t)\}.$$

Take  $\bar{\alpha} := 1 - \alpha$ ,  $I_\alpha$  be an  $n$  by  $n$  diagonal matrix, such that  $(I_\alpha)_{ii} = \bar{\alpha}^{L_i}$ . Let

$$\mathbf{y}_\alpha(t) = \bar{\alpha}^t \mathbf{y}(t) + \mathbf{y}'(t) \quad \text{and} \quad \mathbf{z}_\alpha(t) = \bar{\alpha}^t I_\alpha \mathbf{z}(t) + \mathbf{z}'(t).$$

where  $\{\mathbf{y}'(t)\}$  and  $\{\mathbf{z}'(t)\}$  are specified as follows. For  $t = 1, 2, \dots, L_n - 1$ , we take  $\mathbf{y}'(t) = 0$ . For  $t \geq L_n$ , we take  $\mathbf{y}'(t) := (\bar{\alpha}^{L_n} - \bar{\alpha}^t) \mathbf{d}(t - L_n)$ . The intuition of  $\mathbf{y}_\alpha(t)$  is that the production of  $\bar{\alpha}^t \mathbf{y}(t)$  components can only handle no more than  $\bar{\alpha}^t \mathbf{d}(t)$  out of  $\mathbf{d}(t)$  demands that arrive at period  $t$ . To satisfy the rest  $(1 - \bar{\alpha}^t) \mathbf{d}(t)$  demands, the system schedules a production in period  $t + L_n$  and prepares the required components during period  $[t, t + L_n - 1]$ . Notice that the system only needs to assemble  $\mathbf{y}'(t)$  products due to the presence of damping. The choice of  $\{\mathbf{z}'(t)\}$  is such that the demand arriving at period  $t$  matches the assembly need of  $\mathbf{y}'(t)$ , for all  $t \geq L_n$ . Formally,

$$\mathbf{z}'_i(t) := (A^\top \mathbf{y}'(t + L_i))_i = (\bar{\alpha}^{L_n} - \bar{\alpha}^{t+L_i}) (A^\top \mathbf{d}(t - L_n + L_i))_i, \quad \text{for } t \geq L_n - L_i.$$

Denote  $\mathbf{s}_\alpha(t) = (\mathbf{x}_\alpha(t), \tilde{\mathbf{u}}_\alpha(t), R_\alpha(t))$  for all  $t \geq 0$ . We first show that

$$\mathbf{x}_\alpha(t) = \bar{\alpha}^t \mathbf{x}(t) + A^\top \mathbf{y}'(t) \quad \text{and} \quad \tilde{\mathbf{u}}_\alpha(t) = \bar{\alpha}^t \tilde{\mathbf{u}}(t) + \sum_{s=t}^{t+L_n} \bar{\alpha}^{t-s} \mathbf{y}'(s). \quad (47)$$

Again proof by induction. First, it is clear that the equality holds when  $t = 0$ . Take an induction hypothesis that the equality holds for some  $t_0$  and let  $t = t_0 + 1$ . Expanding  $f_\alpha$ , the choice of  $\mathbf{z}'(t)$  gives that

$$\tilde{\mathbf{u}}_\alpha(t) = \bar{\alpha}(\tilde{\mathbf{u}}_\alpha(t_0) - \mathbf{y}_\alpha(t_0)) + \mathbf{d}(t).$$

In particular, combining the induction hypothesis with the second equation gives

$$\begin{aligned}
\tilde{\mathbf{u}}_\alpha(t) &= \bar{\alpha} \left( \bar{\alpha}^{t-1} \tilde{\mathbf{u}}(t-1) + \sum_{s=t-1}^{t+L_n-1} \bar{\alpha}^{t-1-s} \mathbf{y}'(s) - \mathbf{y}_\alpha(t-1) \right) + \mathbf{d}(t) \\
&= \bar{\alpha}^t (\tilde{\mathbf{u}}(t-1) - \mathbf{y}(t-1) + \mathbf{d}(t)) + \sum_{s=t-1}^{t+L_n-1} \bar{\alpha}^{t-s} \mathbf{y}'(s) - \bar{\alpha} \mathbf{y}'(t-1) + \mathbf{d}(t)(1 - \bar{\alpha}^t) \\
&= \bar{\alpha}^t \tilde{\mathbf{u}}(t) + \sum_{s=t}^{t+L_n-1} \alpha^{t-s} \mathbf{y}'(s) + \bar{\alpha}^{-L_n} \mathbf{y}'(t+L_n) \\
&= \bar{\alpha}^t \tilde{\mathbf{u}}(t) + \sum_{s=t}^{t+L_n} \alpha^{t-s} \mathbf{y}'(s).
\end{aligned}$$

Here, the last equality uses the choice of  $\mathbf{y}'(t-1)$ . Therefore, the second equality of (47) holds for  $t = t_0 + 1$ . For the first equality, the choice of  $\mathbf{z}_\alpha$  results in that the component replenishment received at period  $t$  has two parts. The first part corresponds to the procurement  $\{\bar{\alpha}^t I_\alpha \mathbf{z}(t)\}$ , which thus equals  $\bar{\alpha}^t R(t-1) \mathbf{e}_1$ . The second part corresponds to the procurement  $\{\mathbf{z}'(t)\}$ , which equals  $A^\top \mathbf{y}'(t)$ . As a result, the induction hypothesis again gives

$$\begin{aligned}
\mathbf{x}_\alpha(t) &= \bar{\alpha}(\mathbf{x}_\alpha(t_0) - A^\top \mathbf{y}_\alpha(t_0)) + \bar{\alpha}^t R(t-1) \mathbf{e}_1 + A^\top \mathbf{y}'(t) \\
&= \bar{\alpha}^t (\mathbf{x}(t-1) - A^\top \mathbf{y}(t-1)) + \bar{\alpha}^t R(t-1) \mathbf{e}_1 = \bar{\alpha}^t \mathbf{x}(t) + A^\top \mathbf{y}'(t).
\end{aligned}$$

Therefore, (47) holds for  $t = t_0 + 1$ , and thus for all  $t \geq 0$ .

We next show that  $\pi_\alpha$  is an admissible policy for the damped control problem, namely,  $\pi_\alpha \in \Pi_\alpha$ . for  $\mathbf{a}_\alpha(t) := \{\mathbf{y}_\alpha(t), \mathbf{z}_\alpha(t)\}$ . We need to show that (i)  $(\mathbf{y}_\alpha(t), \mathbf{z}_\alpha(t))$  is independent of  $(\mathbf{d}(t+1), \mathbf{d}(t+2), \dots)$ ; and (ii)

$$\mathbf{a}_\alpha(t) \in \mathcal{A}_\alpha(\mathbf{s}_\alpha(t)), \quad \text{for all } t \geq 0.$$

Here, (i) is clear due to the choice of  $\mathbf{a}_\alpha(t)$  and the admissibility of policy  $\pi = \{\mathbf{a}(t)\}$ . To show (ii), it is known that  $A^\top \mathbf{y}(t) \leq \mathbf{x}(t)$  and  $\mathbf{y}(t) \leq \tilde{\mathbf{u}}(t)$ . Using the choices of  $\mathbf{y}_\alpha(t)$  and (47), it holds that

$$A^\top \mathbf{y}_\alpha(t) \leq \mathbf{x}_\alpha(t) \quad \text{and} \quad \mathbf{y}_\alpha(t) \leq \tilde{\mathbf{u}}_\alpha(t). \quad (48)$$

Therefore,  $\mathbf{a}_\alpha(t) \in \mathcal{A}_\alpha(\mathbf{s}_\alpha(t))$ . It is also clear that  $y_\alpha(t) \geq 0$  and  $z_\alpha(t) \geq 0$  for all  $t \geq 0$ . The choices of  $\{\mathbf{y}'(t)\}$  and  $\{\mathbf{z}'(t)\}$  give  $\|\mathbf{y}'(t)\|_\infty \leq \bar{d}$  and  $\|\mathbf{z}'(t)\| \leq \bar{a}\bar{d}$ . Due to the choice that  $\alpha \leq \bar{d}/\|\mathbf{s}\|_\infty$ ,

$$\|\mathbf{y}(t)\|_\infty \leq \|\tilde{\mathbf{u}}(0)\|_\infty + t\bar{d} \leq (t + \alpha^{-1})\bar{d}.$$

Because  $\pi \in \Pi_S$ , the additional constraint on procurement results in

$$\|\mathbf{z}(t)\|_\infty \leq \bar{a}(\|\tilde{\mathbf{u}}(0)\|_\infty + (t + L_n)\bar{d}) \leq \bar{a}(t + L_n + \alpha^{-1})\bar{d}.$$

We have the inequalities

$$\begin{aligned}
(t + \alpha^{-1})\bar{\alpha}^t &\leq (t + \alpha^{-1})/(1 + \alpha)^t \leq (t + \alpha^{-1})/(1 + \alpha t) \leq \alpha^{-1}, \quad \text{for } t \geq 0, \\
\|\mathbf{y}_\alpha(t)\|_\infty &\leq \bar{\alpha}^t \|\mathbf{y}(t)\|_\infty + \|\mathbf{y}'(t)\|_\infty \leq \bar{\alpha}^t (t + \alpha^{-1})\bar{d} + \bar{d} \leq 2\alpha^{-1}\bar{d} = u_\alpha \quad \text{and} \\
\|\mathbf{z}_\alpha(t)\|_\infty &\leq \bar{\alpha}^t \|\mathbf{z}(t)\|_\infty + \|\mathbf{z}'(t)\|_\infty \leq \bar{a}\bar{\alpha}^t (\alpha^{-1} + t + L_n)\bar{d} + \bar{a}\bar{d} \leq 3\alpha^{-1}L_n\bar{a}\bar{d} = z_\alpha.
\end{aligned}$$

Therefore,  $\mathbf{a}_\alpha(t) \in [0, y_\alpha]^m \times [0, z_\alpha]^n$ , implying that  $\pi_\alpha \in \Pi_\alpha$  is an admissible policy.

Finally, we can expand the expression of  $V^\pi$  and  $V_\alpha^{\pi_\alpha}$  by

$$V^\pi(\mathbf{s}; F) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t C(\mathbf{s}(t), \mathbf{a}(t)) \right] \quad \text{and} \quad V_\alpha^{\pi_\alpha}(\mathbf{s}; F) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t C(\mathbf{s}_\alpha(t), \mathbf{a}_\alpha(t)) \right].$$

We can compute that

$$\begin{aligned}
C(\mathbf{s}_\alpha(t), \mathbf{a}_\alpha(t)) &= \bar{\alpha}^t C(\mathbf{s}(t), \mathbf{a}(t)) + \mathbf{b}^\top \sum_{s=t}^{t+L_n} \alpha^{t-s} \mathbf{y}'(s) \\
&\leq C(\mathbf{s}(t), \mathbf{a}(t)) + \bar{c}\bar{d}m(L_n + 1)(1 - \bar{\alpha}^t).
\end{aligned}$$

Here, the last inequality utilizes the fact that  $\mathbf{y}'(s) = 0$  for  $s < L_n$  and

$$\bar{\alpha}^{t-s} \mathbf{y}'(s) = (\bar{\alpha}^{L_n+t-s} - \bar{\alpha}^t) \mathbf{d}(s - L_n) \leq (1 - \bar{\alpha}^t) \mathbf{1},$$

for  $s \geq L_n$ . Therefore,

$$V_\alpha^{\pi_\alpha}(\mathbf{s}; F) \leq V^\pi(\mathbf{s}; F) + \bar{c}\bar{d}m(L_n + 1) \sum_{t \geq 0} \gamma^t (1 - \bar{\alpha}^t) \leq V^\pi(\mathbf{s}; F) + \frac{\alpha \bar{c}\bar{d}m(L_n + 1)}{(1 - \gamma)^2},$$

giving (46) and finishing the proof.  $\square$

## D.7 The Covering Arguments

We finally prove Lemma 7 using covering arguments. For function  $v : \mathcal{S}_\alpha \rightarrow \mathbb{R}$ , Markovian policy  $\xi$ , and  $\mathbf{s} \in \mathcal{S}$ , let

$$\mathcal{P}^\xi v(\mathbf{s}) := \mathbb{E}_{\mathbf{d}}[v(f_\alpha(\mathbf{s}, \xi(\mathbf{s}), \mathbf{d}))], \quad \text{and} \quad \widehat{\mathcal{P}}^\xi v(\mathbf{s}) := \mathbb{E}_{\hat{\mathbf{d}}}[v(f_\alpha(\mathbf{s}, \xi(\mathbf{s}), \hat{\mathbf{d}}))], \quad (49)$$

Clearly,  $\mathcal{P}^\xi$  and  $\widehat{\mathcal{P}}^\xi$  are state-transition operator with respect to Markovian policy  $\xi$ .

**Lemma 10.** *For Markovian policy  $\xi$  and function  $v : \mathcal{S}_\alpha \rightarrow \mathbb{R}$  such that  $\|v\|_\infty < \infty$  it holds that*

$$\|v\|_\infty \leq (1 - \gamma)^{-1} \|(I - \gamma \mathcal{P}^\xi)v\|_\infty \quad \text{and} \quad \|v\|_\infty \leq (1 - \gamma)^{-1} \|(I - \gamma \widehat{\mathcal{P}}^\xi)v\|_\infty \quad (50)$$

where  $I$  is the identity operator, i.e.,  $Iv = v$  for all  $v$ . Moreover, if  $\pi \in \Pi$  is the induced policy by  $\xi$ , then

$$(I - \gamma \mathcal{P}^\xi)V^\pi(\mathbf{s}; F) = C(\mathbf{s}, \xi(\mathbf{s})) = (I - \gamma \widehat{\mathcal{P}}^\xi)V^\pi(\mathbf{s}; \widehat{F}), \quad \text{for all } \mathbf{s} \in \mathcal{S}_\alpha. \quad (51)$$

*Proof.* Proof of Lemma 10 We only provide the proof for operator  $\mathcal{P}^\xi$ . The results for operator  $\widehat{\mathcal{P}}^\xi$  follow from similar arguments. It follows from (49) that  $\|\mathcal{P}^\xi v\|_\infty \leq \|v\|_\infty$ . Therefore,

$$\|v\|_\infty = \|(I - \gamma \mathcal{P}^\xi)v + \gamma \mathcal{P}^\xi v\|_\infty \leq \|(I - \gamma \mathcal{P}^\xi)v\|_\infty + \gamma \|\mathcal{P}^\xi v\|_\infty \leq \|(I - \gamma \mathcal{P}^\xi)v\|_\infty + \gamma \|v\|_\infty.$$

Rearranging terms,  $\|v\|_\infty \leq (1 - \gamma)^{-1} \|(I - \gamma \mathcal{P}^\xi)v\|_\infty$ , finishing the proof.

Next, we take the induced policy  $\pi$  and the corresponding controlled state process  $\{\mathbf{s}_\alpha(t) : t \geq 0\}$  through (23). It holds that

$$\begin{aligned} V_\alpha^\pi(\mathbf{s}; F) &= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t C(\mathbf{s}_\alpha(t), \xi(\mathbf{s}_\alpha(t))) \middle| \mathbf{s}_\alpha(0) = \mathbf{s} \right] \\ &= C(\mathbf{s}, \xi(\mathbf{s})) + \mathbb{E} \left[ \sum_{t=1}^{\infty} \gamma^t C(\mathbf{s}_\alpha(t), \xi(\mathbf{s}_\alpha(t))) \middle| \mathbf{s}_\alpha(1) = f_\alpha(\mathbf{s}, \xi(\mathbf{s}), \mathbf{d}) \right] \\ &= C(\mathbf{s}, \xi(\mathbf{s})) + \gamma \mathbb{E}_{\mathbf{d}}[V_\alpha^\pi(f_\alpha(\mathbf{s}, \xi(\mathbf{s}), \mathbf{d}); F)] \\ &= C(\mathbf{s}, \xi(\mathbf{s})) + \gamma \mathcal{P}^\xi V_\alpha^\pi(\mathbf{s}; F), \end{aligned}$$

which concludes the proof. □

*Proof.* Proof of Lemma 7. We first prove (30). Rearranging term in (51) gives

$$(I - \gamma \mathcal{P}^\xi)(V_\alpha^\pi(\mathbf{s}; F) - V_\alpha^\pi(\mathbf{s}; \widehat{F})) = \gamma(\mathcal{P}^\xi - \widehat{\mathcal{P}}^\xi)V_\alpha^\pi(\mathbf{s}; \widehat{F}), \quad \text{for all } \mathbf{s} \in \mathcal{S}.$$

Using Lemma 10, we obtain

$$\|V_\alpha^\pi(\cdot; F) - V_\alpha^\pi(\cdot; \widehat{F})\|_\infty \leq \frac{\gamma}{1 - \gamma} \|(\mathcal{P}^\xi - \widehat{\mathcal{P}}^\xi)V_\alpha^\pi(\cdot; \widehat{F})\|_\infty. \quad (52)$$

To bound the right-hand side of (52), it follows from (49) that

$$\begin{aligned} 0 &\leq \mathcal{P}^\xi(V_\alpha^\pi(\mathbf{s}; \widehat{F}) - V_\alpha^*(\mathbf{s}; \widehat{F})) \leq \|V_\alpha^\pi(\mathbf{s}; \widehat{F}) - V_\alpha^*(\mathbf{s}; \widehat{F})\|_\infty, \quad \text{and} \\ 0 &\leq \widehat{\mathcal{P}}^\xi(V_\alpha^\pi(\mathbf{s}; \widehat{F}) - V_\alpha^*(\mathbf{s}; \widehat{F})) \leq \|V_\alpha^\pi(\mathbf{s}; \widehat{F}) - V_\alpha^*(\mathbf{s}; \widehat{F})\|_\infty. \end{aligned}$$

Subtracting the two inequalities and taking supreme over  $\mathbf{s} \in \mathcal{S}_\alpha$ ,

$$\|(\mathcal{P}^\xi - \widehat{\mathcal{P}}^\xi)(V_\alpha^\pi(\cdot; \widehat{F}) - V_\alpha^*(\cdot; \widehat{F}))\|_\infty \leq \|V_\alpha^\pi(\mathbf{s}; \widehat{F}) - V_\alpha^*(\mathbf{s}; \widehat{F})\|_\infty.$$

As a result,

$$\begin{aligned} \|(\mathcal{P}^\xi - \widehat{\mathcal{P}}^\xi)V^\pi(\cdot; \widehat{F})\|_\infty &\leq \|(\mathcal{P}^\xi - \widehat{\mathcal{P}}^\xi)(V^\pi(\cdot; \widehat{F}) - V^*(\cdot; \widehat{F}))\|_\infty + \|\gamma(\mathcal{P}^\xi - \widehat{\mathcal{P}}^\xi)V^*(\cdot; \widehat{F})\|_\infty \\ &= \|V_\alpha^\pi(\mathbf{s}; \widehat{F}) - V_\alpha^*(\mathbf{s}; \widehat{F})\|_\infty + \sup_{g \in \mathcal{G}} |\mathbb{E}_{\mathbf{d}}[g(\mathbf{d})] - \mathbb{E}_{\hat{\mathbf{d}}}[g(\hat{\mathbf{d}})]|, \end{aligned} \quad (53)$$

where the function class  $\mathcal{G}$  is characterizes by all  $g : \mathcal{D} \rightarrow \mathbb{R}$  such that

$$g(\cdot) = \gamma V_\alpha^*(f(\mathbf{s}, \mathbf{a}, \cdot); \widehat{F}) \text{ for some } (\mathbf{s}, \mathbf{a}) \in \mathcal{SA}.$$

Temporarily take  $L_v := (1 - \gamma)^{-1} \bar{c}(\bar{a} + 1)(m + nL_n)$ . It follows from Corollary 1 (taking  $G = \widehat{F}$ ) that  $V^*(\cdot; \widehat{F})$  is  $L_v$ -Lipschitz continuous. As  $f$  is  $(\bar{a} + 1)$ -Lipschitz continuous, every element in  $\mathcal{G}$  is a  $(\bar{a} + 1)L_v$ -Lipschitz continuous function. Take

$$N_0 := \sqrt{2M}(\bar{a} + 1)/\bar{d} \quad \text{and} \quad N = \lceil x_\alpha N_0 + 1 \rceil^{p_S + p_A},$$

where  $\lceil x \rceil$  is the ceiling operator that returns the smallest integer that compares no less than  $x$ , for  $x \in \mathbb{R}$ . Let the finite set  $\{(\mathbf{s}_i, \mathbf{a}_i) : i = 1, 2, \dots, N\}$  form a lattice on  $\mathcal{SA}_\alpha \subseteq [0, x_\alpha]^{p_S+p_A}$  with unit length  $N_0^{-1}$ . In particular,  $\{(\mathbf{s}_i, \mathbf{a}_i) : i = 1, 2, \dots, N\}$  is an  $N_0^{-1}$ -covering of  $\mathcal{SA} \subseteq [0, x_\alpha]^{p_S+p_A}$  with respect to the  $L^\infty$  norm, i.e.,

$$\sup_{(\mathbf{s}, \mathbf{a}) \in \mathcal{SA}} \inf_{i \in [N]} \|(\mathbf{s}, \mathbf{a}) - (\mathbf{s}_i, \mathbf{a}_i)\|_\infty \leq N_0^{-1}.$$

Let  $g_i : \mathcal{D} \rightarrow \mathbb{R}$  such that  $g_i(\mathbf{d}) := V_\alpha^*(f_\alpha(\mathbf{s}_i, \mathbf{a}_i, \mathbf{d}); \widehat{F})$ , for  $i \in [N]$ , the  $(\bar{a} + 1)L_v$ -Lipschitz continuity of the function  $(\mathbf{s}, \mathbf{a}, \mathbf{d}) \mapsto V_\alpha^*(f_\alpha(\mathbf{s}, \mathbf{a}, \mathbf{d}); \widehat{F})$  gives

$$\sup_{g \in \mathcal{G}} \inf_{i \in [N]} \|g - g_i\|_\infty = \sup_{(\mathbf{s}, \mathbf{a}, \mathbf{d}) \in \mathcal{SA} \times \mathcal{D}} \inf_{i \in [N]} |V_\alpha^*(f_\alpha(\mathbf{s}, \mathbf{a}, \mathbf{d}); \widehat{F}) - V_\alpha^*(f_\alpha(\mathbf{s}_i, \mathbf{a}_i, \mathbf{d}); \widehat{F})| \leq (\bar{a} + 1)L_v/N_0.$$

Therefore, for any individual function  $g \in \mathcal{G}$ , there exists an  $i \in [N]$  such that

$$\|g_i - g\|_\infty \leq (\bar{a} + 1)L_v/N_0.$$

In particular,

$$\left| \mathbb{E}_{\mathbf{d}}[g(\mathbf{d}) - g_i(\mathbf{d})] \right| \leq (\bar{a} + 1)L_v/N_0 \quad \text{and} \quad \left| \mathbb{E}_{\hat{\mathbf{d}}}[g_i(\hat{\mathbf{d}}) - g(\hat{\mathbf{d}})] \right| \leq (\bar{a} + 1)L_v/N_0,$$

implying that

$$\begin{aligned} \left| \mathbb{E}_{\mathbf{d}}[g(\mathbf{d})] - \mathbb{E}_{\hat{\mathbf{d}}}[g(\hat{\mathbf{d}})] \right| &\leq \left| \mathbb{E}_{\mathbf{d}}[g(\mathbf{d}) - g_i(\mathbf{d})] \right| + \left| \mathbb{E}_{\mathbf{d}}[g_i(\mathbf{d})] - \mathbb{E}_{\hat{\mathbf{d}}}[g_i(\hat{\mathbf{d}})] \right| + \left| \mathbb{E}_{\hat{\mathbf{d}}}[g_i(\hat{\mathbf{d}}) - g(\hat{\mathbf{d}})] \right| \\ &\leq 2(\bar{a} + 1)L_v/N_0 + \left| \mathbb{E}_{\mathbf{d}}[g_i(\mathbf{d})] - \mathbb{E}_{\hat{\mathbf{d}}}[g_i(\hat{\mathbf{d}})] \right|. \end{aligned}$$

Taking supreme over  $g \in \mathcal{G}$ , we have shown that

$$\sup_{g \in \mathcal{G}} \left| \mathbb{E}_{\mathbf{d}}[g(\mathbf{d})] - \mathbb{E}_{\hat{\mathbf{d}}}[g(\hat{\mathbf{d}})] \right| \leq 2(\bar{a} + 1)L_v/N_0 + \sup_{i \in [N]} \left| \mathbb{E}_{\mathbf{d}}[g_i(\mathbf{d})] - \mathbb{E}_{\hat{\mathbf{d}}}[g_i(\hat{\mathbf{d}})] \right|.$$

Next, we expand the second term on the right-hand side by

$$\mathbb{E}_{\mathbf{d}}[g_i(\mathbf{d})] - \mathbb{E}_{\hat{\mathbf{d}}}[g_i(\hat{\mathbf{d}})] = \sum_{k=1}^M \left( g_i(\mathbf{d}_k) - \mathbb{E}_{\mathbf{d}}[g_i(\mathbf{d})] \right),$$

which, for each  $i \in [N]$ , is the sum of  $M$  copies of mean zero i.i.d. random variables that are bounded by

$$\begin{aligned} &\sup_{\mathbf{d}_1 \in \mathcal{D}} |g_i(\mathbf{d}_1) - \mathbb{E}[g_i(\mathbf{d})]| \\ &= \sup_{\mathbf{d}_1 \in \mathcal{D}} |V_\alpha^*(f_\alpha(\mathbf{s}_i, \mathbf{a}_i, \mathbf{d}_1), \widehat{F}) - \mathbb{E}_{\mathbf{d}}[V_\alpha^*(f_\alpha(\mathbf{s}_i, \mathbf{a}_i, \mathbf{d}), \widehat{F})]| \\ &\leq \sup_{\mathbf{d}_1 \in \mathcal{D}} \mathbb{E}_{\mathbf{d}} \left[ |V_\alpha^*(f_\alpha(\mathbf{s}_i, \mathbf{a}_i, \mathbf{d}_1), \widehat{F}) - V_\alpha^*(f_\alpha(\mathbf{s}_i, \mathbf{a}_i, \mathbf{d}), \widehat{F})| \right] \\ &\leq L_v \sup_{\mathbf{d}_1 \in \mathcal{D}} \mathbb{E}_{\mathbf{d}} \left[ \|f_\alpha(\mathbf{s}_i, \mathbf{a}_i, \mathbf{d}_1) - f_\alpha(\mathbf{s}_i, \mathbf{a}_i, \mathbf{d})\|_\infty \right] \leq L_v \bar{d}. \end{aligned}$$

Temporarily take  $B_g := L_v \bar{d}$ . For each  $i \in [N]$ , it follows from Hoeffding's inequality, (see, for example, Proposition 2.5 and (2.11) in [14]),

$$\mathbb{P} \left( \left| \mathbb{E}_{\mathbf{d}}[g_i(\mathbf{d})] - \mathbb{E}_{\hat{\mathbf{d}}}[g_i(\hat{\mathbf{d}})] \right| \geq \delta_0 \right) \leq \exp \left( -\frac{2M\delta_0^2}{B_g^2} \right),$$

for all  $\delta_0 > 0$ . Here  $\mathbb{P}$  is taken with respect to  $\widehat{D}$ . Taking  $\delta = N \exp(-2B_g^{-2}M\delta_0^2)$ ,

$$\sup_{i \in [N]} \left| \mathbb{E}_{\mathbf{d}}[g_i(\mathbf{d})] - \mathbb{E}_{\hat{\mathbf{d}}}[g_i(\hat{\mathbf{d}})] \right| \leq \delta_0 = M^{-1/2} B_g (\log(N\delta^{-1})/2)^{1/2}$$

with probability at least  $1 - \delta$ . Together with (52) and (53), we have shown that the event

$$\begin{aligned} \left\| V_\alpha^\pi(\cdot; F) - V_\alpha^\pi(\cdot; \widehat{F}) \right\|_\infty &\leq \frac{\gamma}{1-\gamma} \left( \left\| V_\alpha^\pi(\cdot; \widehat{F}) - V_\alpha^*(\cdot; \widehat{F}) \right\|_\infty \right. \\ &\quad \left. + 2(\bar{a} + 1)L_v/N_0 + (2M)^{-1/2} B_g (\log(N\delta^{-1}))^{1/2} \right) \end{aligned} \quad (54)$$

holds with probability at least  $1 - \delta$ .

Finally, the choice of  $N_0$  gives  $(\bar{a} + 1)L_v/N_0 \leq (2M)^{-1/2} B_g$ . It follows from the choice of  $x_\alpha$  that  $x_\alpha \geq \bar{d}$ , and thus  $N_0 x_\alpha \geq \sqrt{2M}(\bar{a} + 1) \geq 2\sqrt{2}$ . For  $\delta \in (0, 1)$ , we have

$$\log(N\delta^{-1}) \geq \log(N_0 x_\alpha + 1) \geq \log(2\sqrt{2} + 1) \geq 1,$$

giving

$$(\bar{a} + 1)L_v/N_0 \leq (2M)^{-1/2} B_g(\log(N\delta^{-1}))^{1/2}.$$

We also have

$$N_0 x_\alpha + 2 \leq x_\alpha(\sqrt{2M}(\bar{a} + 1) + 2)/\bar{d} \leq 3(\bar{a} + 1)M/\bar{d},$$

giving

$$\log(\delta^{-1}N) \leq (p_S + p_A) \log(\delta^{-1}(N_0 x_\alpha + 2)) \leq (p_S + p_A) \log \frac{3(\bar{a} + 1)M x_\alpha}{\delta \bar{d}}$$

Therefore, the “constant” terms in the right-hand side of (54) can be bounded by

$$\begin{aligned} & \frac{\gamma}{1-\gamma} \left( 2(\bar{a} + 1)L_v/N_0 + (2M)^{-1/2} B_g(\log(N\delta^{-1}))^{1/2} \right) \\ & \leq \frac{\gamma}{1-\gamma} \left( 3L_v \bar{d} \sqrt{\frac{p_S + p_A}{2M} \log \frac{3M(\bar{a} + 1)x_\alpha}{\delta \bar{d}}} \right) \\ & \leq \frac{3(\bar{a} + 1)\bar{c}p_S}{(1-\gamma)^2} \sqrt{\frac{p_S + p_A}{2M} \log \frac{3M(\bar{a} + 1)x_\alpha}{\delta \bar{d}}} = \mathcal{E}_1(\delta). \end{aligned}$$

Together with (54), we obtain (52).

The proof of (31) is similar. Using Lemma 10, similar computation to (53) gives

$$\|V_\alpha^\pi(\cdot; F) - V_\alpha^\pi(\cdot; \hat{F})\|_\infty \leq \frac{\gamma}{1-\gamma} \|V_\alpha^\pi(\mathbf{s}; \hat{F}) - V_\alpha^*(\mathbf{s}; \hat{F})\|_\infty + \frac{1}{1-\gamma} \sup_{g \in \mathcal{G}'} |\mathbb{E}_{\mathbf{d}}[g(\mathbf{d})] - \mathbb{E}_{\hat{\mathbf{d}}}[g(\hat{\mathbf{d}})]|,$$

where the class  $\mathcal{G}'$  is characterized by all functions  $g : \mathcal{D} \rightarrow \mathbb{R}$  such that

$$g(\cdot) = \gamma V_\alpha^*(f(\mathbf{s}, \mathbf{a}, \cdot); F) \text{ for some } (\mathbf{s}, \mathbf{a}) \in \mathcal{SA}.$$

Taking  $G = F$  in Corollary 1,  $V_\alpha^*(\cdot; F)$  is  $L_v$ -Lipschitz continuous on  $\mathcal{S}_\alpha$ . Again, similar arguments to the proof of (52) give (31).  $\square$

## E Numerical Results

Hyperparameters	Value
Batch size	$2^8$
Buffer size	$2^{12}$
Soft update	$2^{-8}$
Reward scale	$10^{-2}$
Learning rate	$lr \in [10^{-5}, 10^{-3}]$
Network width	$2^8$
Delayed update	16
Learning Discount	0.95
Number of layers	3

Table 1: Hyperparameter settings for RL algorithms.

**Example 1.** (N-SYSTEM) There are two components and two products in the N-system. Assembling one unit of product 1 requires one unit of component 1. Assembling one unit of product 2 requires one unit of component 1 and one unit of component 2. The inventory holding and backorder costs are  $(h_1, h_2) = (2, 3)$  and  $(b_1, b_2) = (9, 6)$ , respectively. The component lead times are  $L_1 = L_2 = 1$ . Product demands  $d_1$  and  $d_2$  have marginal uniform distributions on  $[0, 40]$ . We separately consider two scenarios where  $d_1$  and  $d_2$  are independent and correlated. When  $d_1$  and  $d_2$  are independent, it is known that the optimal ordering policy follows a base-stock policy with base-stock levels  $s_1 = 51.0$  and  $s_2 = 23.9$  (see [10]). When  $d_1$  and  $d_2$  are correlated, we consider that the conditional distribution of  $[d_2|d_1]$  is  $40 - d_1$  with probability 0.9 and an independent uniform random variable on  $[0, 40]$  with probability 0.1. In particular, the two products become substitutes. Table 2 summarizes the numerical performance of the three policies in different sample sizes.

		Demand Sample Size			
		Algorithm	$M = 10$	$M = 100$	$M = 200$
			$C$ (s.d.)	$C$ (s.d.)	$C$ (s.d.)
I.I.D. Demand (Optimality: 81.0)	CTD3	86.9 (0.04)	82.9 (0.06)	82.2 (0.04)	
	PTD3	87.4 (0.04)	84.6 (0.05)	83.0 (0.04)	
	NV-PRP	86.1 (0.05)	84.7 (0.03)	82.9 (0.04)	
Correlated Demand	CTD3	66.1 (0.03)	60.2 (0.05)	58.1 (0.02)	
	PTD3	67.3 (0.03)	62.6 (0.04)	61.1 (0.05)	
	NV-PRP	79.2 (0.07)	72.8 (0.04)	68.2 (0.08)	

Table 2: The performances of CTD3 policy, PTD3 policy, and NV-PRP policy on the N-system with different demand sample sizes and demand joint distributions.

From Table 2, we first observe that all three methods generate near-optimal policies when the demands of the two products are independent. With  $M = 10$ , the relatively high optimality gap of the two RL algorithms likely results from an overfit of the 10 observed demands. Note that the effect of overfitting is analyzed in [15]. As  $M$  increases, the performance of RL algorithms improves substantially, where the CTD3 algorithm generates a policy with 1.4% (82.2 compared to 81.0) optimality gap. Table 2 also suggests that the CTD3 policy outperforms the PTD3 policy by a statistically significant margin. To this end, Table 2 suggests that RL algorithms are capable of generating near-optimal policies with a relatively small sample of demands.

With independent demand, the near-optimal performance of the NV-PRP policy is expected from its design. Compared to the RL algorithms, the heuristic relies only on the mean demand and, therefore, does not overfit the small sample as much. From Table 2, the NV-PRP policy outperforms the two RL policies by a statistically significant margin when  $M = 10$ . However, due to the suboptimality of the NV-PRP policy, it is eventually outperformed by the two RL policies as  $M$  grows. With correlated demands, the optimal policy is unavailable, and the NV-PRP policy is expected to be costly (see [7]). Indeed, the policies from the two RL algorithms substantially outperform NV-PRP, with 17% cost reduction when  $M = 10$  and around 14% reduction when  $M = 200$ . As a result, our numerical findings suggest that RL algorithms are robust to demand correlation, a common feature that practical ATO systems share.

For a more detailed comparison between the CTD3 algorithm and the PTD3 algorithm, Figure 1 (left and middle) summarizes the convergence rate of the algorithms with different choices of sample size  $M$ . As both CTD3 and PTD3 are iterative algorithms, the comparison is conducted for each iteration count  $K$ . Figure 1 (left and middle) shows that the convergence rate of the CTD3 algorithm is significantly faster than that of the PTD3 algorithm, for both  $M = 10$  and  $M = 200$ , evidencing an improved computational efficiency from incorporating the shape constraint into the neural network architecture.

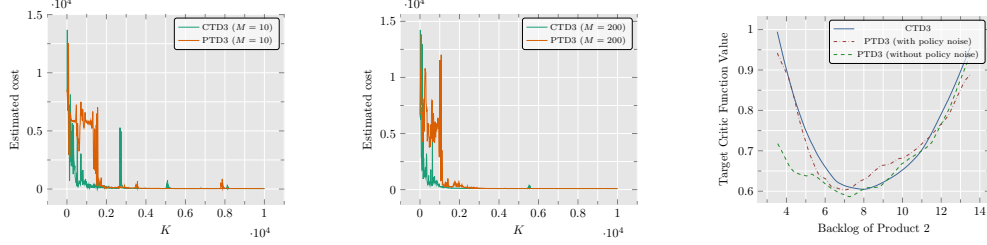


Figure 1: (Color online) Convergence of CTD3 and PTD3 algorithms (left and middle). A slice of actor-value function (right).

The acceleration of convergence can be partly illustrated by Figure 1 (right), which plots a slice of the actor-value function with different  $b_2$ , i.e., the backlog of product 2, while keeping the other state variables fixed. In the figure, the actor-value function computed by the PTD3 algorithm is not convex and is less smooth compared to that of the CTD3 algorithm. Indeed, the policy optimization in the PTD3 algorithm depends on the policy noise, a hyperparameter that smooths the computed value function, so that the optimization is less likely to be stuck in local minima, at the cost of an additional source of error. When policy noise is set to 0, the actor-value function is even less smooth compared to the default value of 0.01. In contrast, the actor-value functions are guaranteed to be convex for the CTD3 algorithm, so that the policy noise can safely be set to 0. To this end, CTD3 delivers a better fit of the actor-value functions compared to PTD3, partially contributing to convergence acceleration.

**Sensitivity Analysis.** Next, we consider how the performance of the two RL algorithms depends on two important tuning factors, the learning rate and the network width. In particular, the learning rate has the same effect on the CTD3 and PTD3 algorithms as the step size affects the convergence of a gradient descent iteration, i.e., a higher learning rate speeds up the initial convergence of the algorithms at the cost of numerical stability. The network width determines the neural network’s size, namely, the number of parameters it contains, with the default network set to 256. For a fixed sample size  $M = 200$  and independent demand, the simulated cost and the optimality gap  $\Delta = [C(\pi) - C(\pi^*)]/C(\pi^*)$  are summarized in Table 3.

Algorithm	Learning Rate	Network Width (NW)					
		NW=128		NW=256		NW=512	
		$C$ (s.d.)	$\Delta$ (%)	$C$ (s.d.)	$\Delta$ (%)	$C$ (s.d.)	$\Delta$ (%)
CTD3	$lr = 5e-5$	84.9 (0.02)	4.60	84.6 (0.03)	4.22	85.1 (0.04)	4.84
	$lr = 1e-4$	82.3 (0.04)	1.53	82.9 (0.03)	2.27	83.1 (0.02)	2.53
	$lr = 5e-4$	84.8 (0.02)	4.51	82.7 (0.05)	2.07	82.7 (0.05)	2.06
	$lr = 1e-3$	82.9 (0.03)	2.34	82.2 (0.04)	1.44	82.2 (0.03)	1.45
PTD3	$lr = 5e-5$	86.2 (0.03)	6.01	86.5 (0.03)	6.32	85.9 (0.05)	5.70
	$lr = 1e-4$	85.2 (0.02)	4.92	84.9 (0.03)	4.62	83.8 (0.04)	3.39
	$lr = 5e-4$	83.7 (0.05)	3.27	83.0 (0.04)	2.41	83.1 (0.03)	2.53
	$lr = 1e-3$	83.3 (0.02)	2.76	84.0 (0.04)	3.55	85.5 (0.04)	5.22

Table 3: Hyperparameter tuning for CTD3 and PTD3 algorithms.

From Table 3, we conclude that the performances of the CTD3 and PTD3 algorithms do not significantly depend on the network width, where both algorithms output near-optimal policies. Hence, in subsequent experiments, we keep this parameter at its default value. On the other hand, the performance of the PTD3 policy depends on the choice of learning rate when the network width is 128 or 256, while the same dependence of the CTD3 policy is not significant.

**Example 2.** (PC-SYSTEM) This example is a simplified PC system with six components and four products considered in [11]. The demand for product  $i \in \{1, 2, 3, 4\}$  follows an independent Poisson distribution of intensity  $\lambda_i$ . We present the BOM matrix and related parameters of this medium-scale PC system in the appendix.

We start with the comparison of the cost among policies obtained from the CTD3 algorithm, the PTD3 algorithm, and the NV-PRP policy, and summarize the results for various sample sizes  $M$  by Figure 2. We first observe that the cost of the CTD3 policy again decreases as  $M$  increases from 10 to 1000. The decrease rate is much faster when  $M \leq 200$ , supporting our result in Theorem 3 that the effective learning of the optimal policy  $\pi^*$  does not require a very large demand sample. Second, Figure 2 (left) indicates that the CTD3 policy consistently achieves a substantially lower cost among all sample sizes compared to the PTD3 policy and the NV-PRP policy for any choice of sample size. For example, when  $M = 100$ , the cost reduction of the CTD3 policy is 3% compared to the PTD3 policy and 9% compared to the NV-PRP policy.

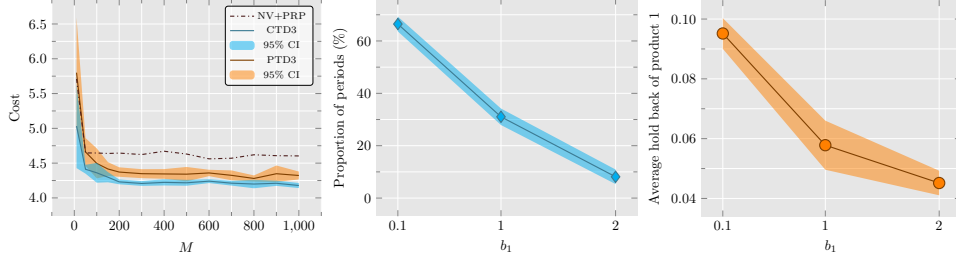


Figure 2: (Color online) Performances of CTD3 and PTD3 algorithms (left). Strategic holdback of product 1 (middle and right).

To analyze the effects of hyperparameter tuning, Table 4 presents the average and standard deviation of the simulated cost for both the PTD3 and CTD3 algorithms under various learning rates from  $\{1e-5, 5e-5, 1e-4, 5e-4\}$  and sample sizes  $M \in \{10, 500, 1000\}$ . We find that for all sample sizes, the CTD3 algorithm outperforms the PTD3 algorithm. Furthermore, the CTD3 algorithm converges for all learning rates, whereas the PTD3 algorithm fails to converge at  $lr = 1e-5$ . The findings, therefore, support our insight that geometric properties of the ATO control problem help improve the performance, as well as the robustness of the general-purpose PTD3 algorithm.

Algorithm	Learning Rate	Demand Sample Size		
		$M = 10$	$M = 500$	$M = 1000$
		$C$ (s.d.)	$C$ (s.d.)	$C$ (s.d.)
CTD3	$lr = 1e-5$	5.78 (0.03)	4.22 (0.03)	4.35 (0.02)
	$lr = 5e-4$	5.03 (0.04)	4.36 (0.02)	4.18 (0.03)
	$lr = 1e-4$	5.26 (0.03)	4.52 (0.03)	4.33 (0.03)
	$lr = 5e-4$	5.43 (0.03)	4.28 (0.03)	4.29 (0.03)
PTD3	$lr = 1e-5$	-*	-	-
	$lr = 5e-5$	5.93 (0.05)	4.90 (0.05)	4.32 (0.03)
	$lr = 1e-4$	5.80 (0.03)	4.34 (0.03)	4.81 (0.03)
	$lr = 5e-4$	6.15 (0.05)	5.47 (0.06)	4.76 (0.04)

\* -: Algorithm fails to converge.

Table 4: Learning-rate tuning of PC-system.

**Performance Analysis of the No-holdback Rule.** The output policy of the highly efficient CTD3 algorithm facilitates a numerical analysis of the suboptimality of the no-holdback rule. During the 10,000-period simulation of the CTD3 policy with sample size  $M = 1000$ , the decisions intriguingly made in 812 periods violate the no-holdback rule by withholding the production of product 1. Indeed, by strategically withholding components required to assemble product 1, the system can save components 2 and 5 for future assembly of more expensive products. To illustrate this effect, we take  $b_1 \in \{0.1, 1, 2\}$  while keeping the other parameters unchanged. Figure 2 (middle and right) summarizes the holdback of product 1, which is averaged over 100 simulation runs, accompanied by a 95% confidence interval by the shaded areas. In particular, the proportion of periods with holdback increased to 31.04% when  $b_1$  decreased from 2 to 1 and to 66.46% when  $b_1 = 0.1$ . To this end, our numerical results support the observation that, when holding costs are relatively low, policies that involve strategic holdback may be more cost-effective than no-holdback policies.

**Example 3.** (LARGE SYSTEM) We finally analyze the performance of the CTD3 algorithm on the ATO system from the processed food company. There are 15 components and 8 products, where Table 5 provides the BOM matrix, lead times, and other cost parameter settings. As discussed, the actual demand data in this example contains 100 monthly demand points for each end product. As a result, the performance metric  $C(\pi)$  is unknown. We report the performance of the PTD3 algorithm, CTD3 algorithms, and the NV-PRP policy by  $\hat{C}(\pi)$ .

From Table 6, the CTD3 policy (tuned) yields an average cost of 32.0, outperforming the NV-PRP policy by 11.6%. We also observe that the CTD3 algorithm can achieve satisfying policies under the learning rates from  $\{5e-5, 1e-4, 5e-4\}$ , while the PTD3 algorithm fails to deliver a substantial cost reduction under all choices of the learning rate. It is essential to reiterate that the only difference between PTD3 and CTD3 lies in the network architecture, ICNN versus a general-purpose dense neural network. To this end, our numerical results suggest that the utility of ICNN, which originates from the convex structure of the ATO transition dynamics, becomes more evident as the size of the ATO control problem increases. In the  $N$ -system, the effect of ICNN is limited to convergence rate acceleration. For the slightly larger PC system, the ICNN starts to deliver better policies

Component	$h_j$	$L_j$	Product							
			1	2	3	4	5	6	7	8
1	0.163	2	1	1	1	0	0	0	0	0
2	0.318	1	0	0	0	1	1	1	0	0
3	0.104	1	0	0	0	0	0	0	1	0
4	0.206	2	0	0	0	0	0	0	0	1
5	0.167	1	1	1	0	0	0	0	0	0
6	0.575	1	0	0	1	1	1	0	0	0
7	0.126	1	0	0	0	0	0	1	1	0
8	0.166	1	0	0	0	0	0	0	0	1
9	0.275	1	1	0	1	0	0	0	0	0
10	0.303	2	0	1	0	0	0	0	0	0
11	0.191	1	0	0	0	1	0	1	0	0
12	0.307	1	0	0	0	0	0	0	1	1
13	0.473	1	0	0	0	0	1	0	0	0
14	0.239	1	1	1	0	0	0	0	0	0
15	0.129	3	1	0	0	1	0	1	0	0

Table 5: Large scale PC example.

compared to dense neural networks. The ICNN eventually becomes the critical factor for an RL algorithm to generate cost-efficient policies that outperform the benchmark policies.

	PTD3	CTD3	NV-PRP
	$\bar{C}$ (s.d.)	$\bar{C}$ (s.d.)	$\bar{C}$ (s.d.)
$lr = 1e-5$	-	50.0 (0.16)	
$lr = 5e-5$	-	32.1 (0.05)	
$lr = 1e-4$	77.7 (0.09)	32.0 (0.04)	36.2 (0.14)
$lr = 5e-4$	118.7 (0.07)	33.2 (0.06)	

Table 6: Learning-rate tuning of processed food system.

Figure 3 compares the convergence of CTD3 and PTD3 under different learning rates. We observe that, as the learning rate becomes larger, the initial convergence of the CTD3 algorithm becomes faster but less stable. However, the learning curve of the PTD3 algorithm oscillates in the order of  $10^6$  throughout the iteration process under every choice of the learning rate, where the output policy significantly underperforms that of the CTD3 algorithm. We remark that training neural networks is demanding in computational resources, emphasizing the value of effective algorithmic design. To this end, our numerical results underscore the significance of integrating structural insights into RL algorithms, as seen in architectural variations like ICNN, which becomes increasingly valuable when confronted with intricate real-world control problems.

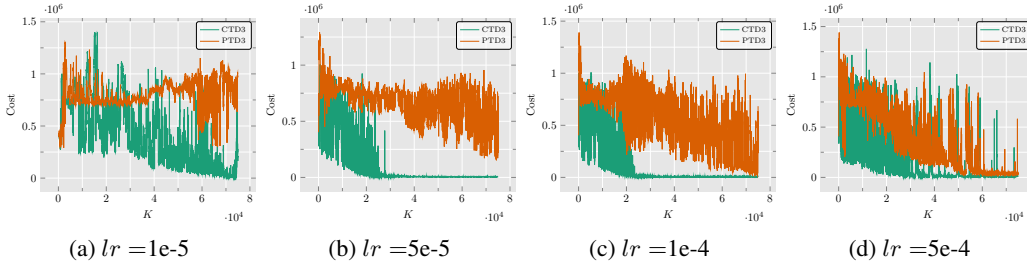


Figure 3: (Color online) Convergence of CTD3 and PTD3 algorithms in the processed food system.