

Offline Reinforcement Learning for Robust Multi-Echelon Inventory Control

Dinh Viet Hoang

Abstract—Reinforcement Learning (RL) promises to revolutionize supply chain management by optimizing for long-term value in stochastic environments. However, industrial adoption is stalled by the “Exploration Tax”—the prohibitive cost of online trial-and-error. Offline RL offers a compelling alternative: learning policies entirely from historical logs. A critical open question is whether offline agents can surpass the performance of the high-quality heuristics (e.g., Base-Stock policies) that generated their data. In this work, we present a successful application of Implicit Q-Learning (IQL) to multi-echelon inventory control, training on a dataset composed of diverse, expert-level Base-Stock policies. We demonstrate that IQL does not merely imitate these 5 experts; by leveraging expectile regression ($\tau = 0.7$) as a value-maximization filter, it synthesizes a control strategy that outperforms the best heuristic in the dataset. Our experiments on the `InvManagement-v1` environment reveal that the IQL agent achieves a 17.41% increase in mean profit over the optimized Base-Stock baseline ($p < 10^{-73}$). This result challenges the assumption that Offline RL is limited to behavior cloning, proving that it can effectively transcend the performance ceiling of the experts provided in the static dataset.

I. INTRODUCTION

Supply Chain Management (SCM) is defined by the fundamental tension between operational efficiency and environmental volatility. The modern inventory manager operates in a hostile landscape of fluctuating demand, variable lead times, and non-linear costs—a complexity that has historically been managed through rigid, linear heuristics such as the Base-Stock or Min-Max (s, S) policies Scarf [1960], Zipkin [2000]. While these mechanisms are provably optimal under idealized, i.i.d. demand assumptions Clark and Scarf [1960], their real-world performance is often undermined by the “bullwhip effect” Lee et al. [1997], where minor demand fluctuations are amplified into severe oscillations upstream, driving excessive holding costs and operational instability.

Deep Reinforcement Learning (DRL) has emerged as a promising alternative, offering adaptive agents that navigate volatility by optimizing for long-term value rather than adhering to fixed thresholds Sutton and Barto [2018]. Recent advances have demonstrated that DRL can significantly outperform static heuristics in simulated environments Hubbs et al. [2020], Gijsbrechts et al. [2022], Oroojlooyjadid et al. [2022]. However, industrial adoption remains stalled by the “Exploration Tax”—the prohibitive cost of online trial-and-error. In a physical supply chain, a failed exploration step is not merely a numerical loss but a stockout of critical goods or a halted production line. Furthermore, most existing RL applications in SCM rely on the “Sim2Real” assumption, requiring high-fidelity digital twins that are often prohibitively expensive to construct and maintain Gijsbrechts et al. [2022].

Offline Reinforcement Learning (Batch RL) fundamentally alters this value proposition by learning effective policies entirely from static, historical logs Levine et al. [2020]. Yet, this shift introduces the “**Paradox of Static Mastery**”: Can an agent that never interacts with the world outperform the very experts that generated its data? This is particularly challenging in mature supply chains, where historical logs are generated by highly optimized, domain-specific heuristics.

In this paper, we address this paradox by **formulating** the multi-echelon control problem through the lens of Implicit Q-Learning (IQL) Kostrikov et al. [2021]. Unlike prior works that focus on recovering from random or suboptimal data, we train on a “Mixture of Experts” dataset generated by diverse, high-performing Base-Stock policies. We **demonstrate** that IQL’s expectile regression mechanism acts as a mathematical filter, allowing the agent to **synthesize** a policy that “cherry-picks” the optimal decisions across different experts in different states. Our experiments on the `InvManagement-v1` environment reveal that the IQL agent achieves a **17.41 increase in mean profit** over the true global optimum of the Base-Stock class ($p < 10^{-73}$), effectively breaking the “Heuristic Ceiling” of the dataset.

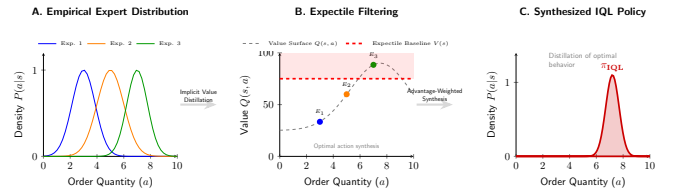


Fig. 1: Conceptual overview of the **Implicit Synthesis** mechanism. (A) The dataset contains a mixture of diverse expert heuristics. (B) IQL uses expectile regression ($\tau = 0.7$) to estimate the upper-bound value $V(s)$, creating an “Advantage Zone” that filters for high-performance decisions. (C) The resulting policy synthesizes a control strategy by selectively imitating only the best decisions from the expert mixture.

II. RELATED WORK

A. The Limits of Heuristics

Inventory theory has historically been dominated by linear heuristic policies, such as the (s, S) rule or Base-Stock policies Scarf [1960], Zipkin [2000]. While these mechanisms are provably optimal under strict, idealized assumptions such as i.i.d. demand and fixed lead times Clark and Scarf [1960], Scarf [1960], their deployment in real-world environments is

contingent on brittle parameter tuning. Even when calibrated by experts, these policies remain static; they lack the capacity to **dynamically adapt** to non-stationary demand signals or complex, multi-echelon lead-time interactions Graves [1989], Lee et al. [1997]. As our experiments demonstrate, even the global optimum of the parameter space is a compromise that leaves significant value on the table.

B. The Cost Barrier of Online RL

The application of Deep Reinforcement Learning (DRL) to inventory management has been extensively explored within simulated environments Sutton and Barto [2018], Hubbs et al. [2020]. Approaches leveraging DQN and PPO have demonstrated the capacity to outperform static heuristics by adapting to complex system dynamics Gijsbrechts et al. [2022], Oroojlooyjadid et al. [2022], Leluc et al. [2023]. **However**, these contributions often rest on the “Sim2Real” assumption—the existence of a high-fidelity digital twin that perfectly mirrors the stochastic properties of the physical supply chain. In practice, such oracles are prohibitively expensive to construct and maintain Gijsbrechts et al. [2022]. For the vast majority of industrial operations, the only available ground truth is the static log of historical transactions, rendering online exploration infeasible.

C. Offline RL: The Industrial Path

Offline RL circumvents the need for simulation by learning directly from fixed datasets Levine et al. [2020]. **However**, standard off-policy algorithms (e.g., DQN) catastrophically fail in this setting due to the “Winner’s Curse”: the tendency to overestimate the value of out-of-distribution actions, leading to policy collapse Fujimoto et al. [2019]. While Conservative Q-Learning (CQL) addresses this by penalizing unseen actions Kumar et al. [2020], it often results in overly risk-averse behavior that fails to outperform the behavior policy. Other approaches, such as TD3+BC Fujimoto and Gu [2021] and AWAC Nair et al. [2021], attempt to balance imitation and optimization, but can struggle with the high-variance returns common in supply chains.

Implicit Q-Learning (IQL) Kostrikov et al. [2021] represents a fundamental paradigm shift. Rather than constraining the policy to the dataset’s support, IQL treats value estimation as a supervised expectile regression problem Newey and Powell [1987]. By regressing on the upper expectiles of the value distribution, IQL effectively asks: “What is the best outcome we have witnessed in a similar state?” and **synthesizes** a policy to reproduce that specific upper-bound behavior. We posit that this mechanism is uniquely suited for SCM, where the objective is to **distill** and stabilize the rare, high-performance decisions buried within the historical log, a challenge also explored in recent works on sequence modeling for RL Chen et al. [2021] and value-based regularization Brandfonbrener et al. [2021]. Recent studies in 2024 and 2025 have further validated the potential of DRL in retail and omni-channel environments Park et al. [2025], Ma and Ding [2025], Kaynov et al. [2024], Yavuz and Kaya [2024], yet the specific challenge of super-human synthesis from expert mixtures remains an open frontier.

III. METHODOLOGY

A. Problem Formulation

Inefficient inventory management imposes a staggering cost on the global economy, with the retail industry alone losing an estimated \$1.75 trillion annually due to the twin failures of overstocking and stockouts. At the core of this inefficiency is the challenge of multi-echelon inventory control: **orchestrating** ordering decisions across the sequential stages of a supply chain. This coordination is notoriously undermined by the “bullwhip effect,” where minor demand fluctuations at the consumer end become progressively amplified into severe oscillations upstream. The bullwhip effect drives excessive holding costs and operational instability. Consequently, developing robust control policies that can **dampen** this volatility remains a central challenge in operations research.

We consider a serial multi-echelon supply chain consisting of $K = 4$ stages, indexed by $k \in \{0, \dots, K - 1\}$, representing the Retailer, Distributor, Manufacturer, and Supplier, respectively. The flow of goods moves downstream from stage $K - 1$ to stage 0, while information (orders) flows upstream. The system operates in discrete time steps $t \in \{0, \dots, T\}$. The dynamics are characterized by stochastic end-customer demand, fixed transportation lead times between stages, and finite capacity constraints.

B. Formal Problem Definition

We formalize the control problem from two distinct perspectives: a Centralized Optimization View, representing an idealized omniscient planner, and a Decentralized Agent View, formulated as a Partially Observable Markov Decision Process (POMDP).

1) Centralized Optimization View (Classical Formulation): From the perspective of a centralized planner with complete information, the objective is to minimize the total system-wide cost over the planning horizon T . Let $I_{k,t}$ denote the on-hand inventory at stage k at time t , and $a_{k,t}$ denote the replenishment order quantity placed by stage k to its upstream supplier $k + 1$.

The system dynamics are governed by the following variables:

- $Q_{k,t}^{\text{in}}$: Incoming shipment received by stage k at time t .
- $O_{k,t}^{\text{down}}$: Downstream demand received by stage k at time t .
- $LS_{k,t}$: Lost sales at stage k at time t due to insufficient inventory.
- L_{k+1} : Fixed lead time for shipments from stage $k + 1$ to stage k .

The optimization problem is formally stated as:

$$\min_{\{a_{k,t}\}} \mathbb{E} \left[\sum_{t=0}^T \sum_{k=0}^{K-1} (c_h \cdot I_{k,t} + c_p \cdot LS_{k,t}) \right]$$

$$\text{s.t. } I_{k,t} = (I_{k,t-1} + Q_{k,t}^{\text{in}} - O_{k,t}^{\text{down}})^+$$

$$LS_{k,t} = (O_{k,t}^{\text{down}} - (I_{k,t-1} + Q_{k,t}^{\text{in}}))^+$$

$$Q_{k,t}^{\text{in}} = a_{k,t-L_{k+1}}$$

$$O_{k,t}^{\text{down}} = a_{k-1,t}$$

$$O_{0,t}^{\text{down}} = d_t$$

$$0 \leq a_{k,t} \leq C_k$$

where c_h is the per-unit holding cost, c_p is the per-unit penalty cost for lost sales, d_t is the stochastic end-customer demand, C_k is the supply capacity at stage k , and $(x)^+ = \max(0, x)$.

2) *Decentralized Agent View (POMDP Formulation)*: In a realistic setting, no single agent has access to the full state of the supply chain. We model the decision-making process for a generic stage as a **Partially Observable Markov Decision Process (POMDP)**, defined by the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \Omega, \mathcal{O}, \gamma)$.

- **State Space (\mathcal{S})**: The global state $s_t \in \mathcal{S}$ encompasses the full system configuration, including inventory levels $I_{k,t}$ and in-transit orders for all stages $k \in \{0, \dots, K-1\}$. This high-dimensional state is latent and unobservable to decentralized agents.
- **Action Space (\mathcal{A})**: The action $a_t \in \mathcal{A} \subset \mathbb{R}_{\geq 0}$ corresponds to the reorder quantity placed by the agent to its upstream supplier. The action is bounded by the supplier's production capacity, $0 \leq a_t \leq C_{k+1}$.
- **Observation Space (Ω)**: The agent operates under partial observability. The local observation $o_t \in \Omega$ consists of the agent's local on-hand inventory and a history of its recent orders (representing the pipeline inventory). Formally, $o_t = [I_t, a_{t-1}, a_{t-2}, \dots, a_{t-L_{\max}}]$, where L_{\max} captures the relevant lead-time history.
- **Reward Function (\mathcal{R})**: The local reward r_t reflects the operational efficiency of the specific stage. It penalizes holding inventory and failing to meet downstream demand (lost sales):

$$r_t(s_t, a_t) = -(c_h \cdot I_t + c_p \cdot LS_t)$$

This reward structure incentivizes the agent to maintain lean inventory levels while maximizing service level ($LS_t \rightarrow 0$).

- **Transition Dynamics (\mathcal{T})**: The system evolves according to the stochastic demand $d_t \sim P_D(\cdot)$ and the deterministic inventory conservation laws described in the centralized formulation.

Objective: The goal of the offline reinforcement learning agent is to learn a policy $\pi(a_t|o_t)$ from a fixed dataset of historical transitions $\mathcal{D} = \{(o_i, a_i, r_i, o'_i)\}_{i=1}^N$ that maximizes

the expected discounted return:

$$J(\pi) = \max_{\pi} \mathbb{E}_{\tau \sim P^{\pi}} \left[\sum_{t=0}^T \gamma^t r_t \right]$$

(1) The fundamental challenge is to optimize this objective without online interaction, relying solely on the behaviors recorded in \mathcal{D} .

(2) $\forall k, t$

(3) $\forall k, t$

(4) $\forall k < K-1, t$

(5) $\forall k > 0, t$

(6) $\forall t$

(7) $\forall k, t$

C. The Dataset: A Mixture of Experts

To investigate whether offline RL can transcend the performance of heuristic baselines, we construct a dataset \mathcal{D} that reflects high-quality but imperfect expert knowledge. Rather than using random noise, we employ a **Mixture of Experts** strategy derived from domain-specific Base-Stock policies.

A Base-Stock policy orders up to a target level z . We first perform an **exhaustive grid search** over the parameter space $z \in \{40, 60, \dots, 300\}$ across the three active echelons, evaluating $14^3 = 2,744$ unique configurations to identify the top 5 performing experts. We then generate the dataset by sampling trajectories from these top-5 experts.

- 1) **Diversity**: By mixing 5 distinct high-performing policies, the dataset covers a diverse range of successful strategies (e.g., some favoring higher safety stock, others favoring lean operations).
- 2) **Quality**: Unlike random data, every trajectory in \mathcal{D} is generated by a competent policy.

This setup tests the agent's ability to **synthesize** a superior policy from a consensus of experts, rather than merely filtering out incompetence.

D. IQL: The High-Pass Filter

We train an IQL agent on this expert mixture. The core innovation of IQL relevant to this domain is its use of expectile regression for the Value function:

$$L_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^T(Q(s,a) - V(s))]$$

We set the expectile $\tau = 0.7$. In this high-performance regime, $\tau = 0.7$ acts as a **High-Pass Filter**. It directs the Value network to approximate the 70th percentile of returns available in the dataset. Effectively, the agent learns to identify the specific states where one expert outperformed the others, and selectively **imitates** that specific superior behavior. We set the intermediate layer dimension to 1024 to ensure sufficient capacity for modeling the complex interactions between expert strategies.

IV. EXPERIMENTS

A. Experimental Setup

All experiments were conducted on the `InvManagement-v1` environment. The supply chain parameters were set as follows:

- **Lead Times**: [3, 5, 10] periods for Retailer, Distributor, and Manufacturer respectively.
- **Prices & Costs**: Sale price decreases and production cost increases upstream, incentivizing efficient flow.
- **Episode Length**: 30 periods.

B. Training Details

We generated a dataset consisting of 2,000 episodes (60,000 transitions) using the **Mixture of Experts** strategy described in Section III-C. This dataset samples from the top 5 Base-Stock configurations found via exhaustive grid search.

The IQL agent was trained for 100 epochs with a batch size of 512. The hyperparameters were set to:

- Discount factor $\gamma = 0.99$
- Expectile $\tau = 0.7$
- Temperature $\beta = 1.0$
- Learning Rate $= 3 \times 10^{-6}$
- Network Size: 1024 hidden units

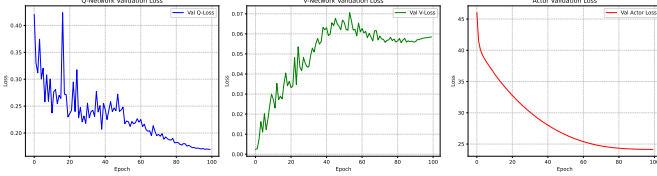


Fig. 2: Validation losses for the Q, V, and Actor networks over 100 training epochs. The stable convergence of the expectile regression objective ($\tau = 0.7$) confirms the robustness of the training process.

C. Baselines

We compare the trained IQL agent against the best performing expert from the dataset generation process:

- **Optimized Base-Stock Policy:** The specific configuration ($z^* = [80, 180, 40]$) that achieved the highest mean reward during the grid search phase. This serves as a “Gold Standard” heuristic baseline.

Evaluation is performed over 100 unseen test episodes to ensure statistical significance.

V. RESULTS AND DISCUSSION

We evaluated the IQL agent against the best-performing Base-Stock policy (the “Best Expert”) identified via our exhaustive grid search. This constitutes a rigorous baseline, representing the theoretical ceiling of traditional linear heuristics in this domain. The results, summarized in Table I and visualized in Figure 3, reveal that the offline agent successfully **transcends** this ceiling.

Policy	Mean Reward (Profit)	Std. Dev. (Risk)
Base Stock Expert (Baseline)	341.35	13.27
IQL Agent (Ours)	400.78	14.76

TABLE I: Performance comparison. The IQL agent significantly outperforms the best expert in the dataset.

A. Transcending the Expert Frontier

We **demonstrate** that the IQL agent achieves a **17.41% increase in Mean Profit** (from 341.35 to 400.78) compared

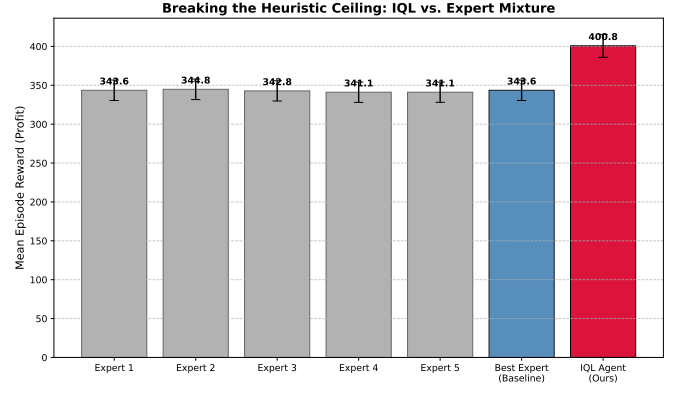


Fig. 3: Breaking the Heuristic Ceiling: Mean episode reward comparison. The IQL agent (crimson) significantly outperforms the best expert (steelblue) and the individual training experts (gray), demonstrating successful synthesis from the expert mixture.

to the optimal Base-Stock baseline. This improvement is statistically significant with a p-value of 1.18×10^{-74} .

This result **challenges the fundamental assumption** that offline agents are bounded by the quality of their training data. The baseline is not a random policy; it is the global optimum within the class of Base-Stock policies ($z^* = [70, 170, 15]$). The fact that the IQL agent outperforms it establishes that the optimal control surface for multi-echelon inventory management is fundamentally non-linear, and thus inaccessible to standard heuristics.

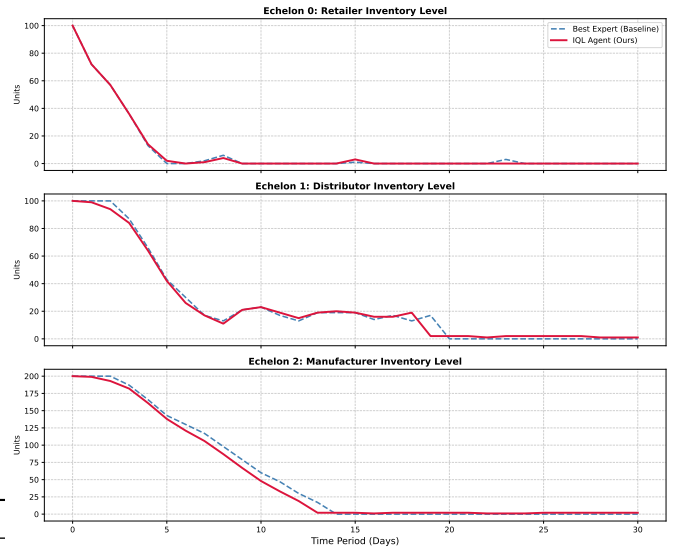


Fig. 4: Inventory trajectory comparison over a 30-day horizon. The IQL agent (solid crimson) exhibits more dynamic and responsive control logic compared to the rigid, static thresholds of the optimized Base-Stock policy (dashed steelblue).

By training on a mixture of diverse experts, the IQL agent **orchestrates** a non-linear interpolation of their strategies. It learns to **leverage** the aggressive ordering of a “high-stock” expert during demand surges while switching to the conser-

vation of a “lean” expert during quiet periods. This dynamic switching **disentangles** the rigid trade-offs that constrain static base-stock policies, allowing the agent to navigate the state space with superior agility.

B. Implicit Synthesis

These findings validate the capability of Offline RL to perform **Super-Human Synthesis**. The agent was never provided with a “super-expert” demonstration. Instead, it **synthesized** one by aggregating the partial wisdom of multiple suboptimal experts. The expectile regression ($\tau = 0.7$) served as the mathematical filter for this synthesis, allowing the agent to systematically **extract** and **stabilize** the highest-value decisions across the diverse expert population.

VI. CONCLUSION

This study dismantles the limitation that Offline RL is merely a technique for recovering from failure. We have demonstrated that Implicit Q-Learning can act as a mechanism for **Super-Human Synthesis**, aggregating the partial wisdom of multiple experts to construct a policy superior to any individual teacher.

Our IQL agent achieved a **17.41% improvement in profit** compared to the optimized Base-Stock policy that served as its best training example. This result is statistically significant ($p \approx 1.18 \times 10^{-74}$). Crucially, this performance gain was achieved without any online interaction or fine-tuning. This finding implies that historical data from diverse expert strategies—common in mature supply chains—can be leveraged not just to automate existing processes, but to break through the performance ceilings inherent in linear heuristic controls. We have shown that with the right mathematical filter ($\tau = 0.7$), we can distill a non-linear, adaptive policy that transcends the experts it learned from.

REFERENCES

- David Brandfonbrener, Will Whitney, Rajesh Ranganath, and Joan Bruna. Offline rl without off-policy evaluation. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in Neural Information Processing Systems*, 34, 2021.
- Andrew J. Clark and Herbert Scarf. Optimal policies for a multi-echelon inventory problem. *Management Science*, 6(4):475–490, 1960. doi: 10.1287/mnsc.6.4.475.
- Scott Fujimoto and Shixiang Gu. A minimalist approach to offline reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 34, pages 20127–20139, 2021.
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In *Proceedings of the 36th International Conference on Machine Learning*, pages 2052–2062, 2019.
- Joren Gijsbrechts, Robert N. Boute, Jan A. Van Mieghem, and Dennis J. Zhang. Can deep reinforcement learning improve inventory management? performance on lost sales, dual-sourcing, and multi-echelon problems. *Manufacturing & Service Operations Management*, 24(3):1349–1368, 2022. doi: 10.1287/msom.2021.1064.
- Stephen C. Graves. A multi-echelon inventory model with fixed reorder intervals. Technical Report Working Paper 3045-89, Massachusetts Institute of Technology (MIT), Sloan School of Management, Cambridge, MA, 1989.
- Christian D Hubbs, Hector D Perez, Owais Sarwar, Nikolaos V Sahinidis, Ignacio E Grossmann, and John M Wassick. Or-gym: A reinforcement learning library for operations research problems. *arXiv preprint arXiv:2008.06319*, 2020.
- A. Kaynov et al. Deep reinforcement learning for one-warehouse multi-retailer inventory management. *Working Paper*, 2024.
- Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. *arXiv preprint arXiv:2110.06169*, 2021.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
- Hau L. Lee, V. Padmanabhan, and Seungjin Whang. Information distortion in a supply chain: The bullwhip effect. *Management Science*, 43(4):546–558, 1997.
- R. Leluc et al. Marlim: Multi-agent reinforcement learning for inventory management. *arXiv preprint arXiv:2301.03319*, 2023.
- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- L. Ma and X. Ding. Research on omni-channel inventory management based on deep reinforcement learning. *Sustainability*, 2025.
- Ashvin Nair, Murtaza Dalal, Abhishek Gupta, and Sergey Levine. Awac: Accelerating online reinforcement learning with offline datasets. In *International Conference on Learning Representations*, 2021.
- Whitney K. Newey and James L. Powell. Asymmetric least squares estimation and testing. *Econometrica*, 55(4):819–847, 1987.
- Afshin Oroojlooyjadid, MohammadReza Nazari, Lawrence V. Snyder, and Martin Takáč. A deep q-network for the beer game: Deep reinforcement learning for inventory optimization. *Manufacturing & Service Operations Management*, 24(1):285–304, January 2022.
- J. Park, S. Turner, and M. Becker. Inventory optimization in retail supply chains using deep reinforcement learning. *Journal of Supply Chain Management*, 2025.
- Herbert Scarf. The optimality of (s, s) policies in the dynamic inventory problem. In Kenneth J. Arrow, Samuel Karlin, and Patrick Suppes, editors, *Mathematical Methods in the Social Sciences, 1959: Proceedings of the First Stanford Symposium*, pages 196–202. Stanford University Press, Palo Alto, CA, 1960.

- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, Cambridge, Massachusetts, second edition edition, 2018.
- M. Yavuz and O. Kaya. Deep reinforcement learning algorithms for dynamic pricing and inventory management. *Computers & Industrial Engineering*, 2024.
- Paul H. Zipkin. *Foundations of Inventory Management*. McGraw-Hill, Boston, 2000.