

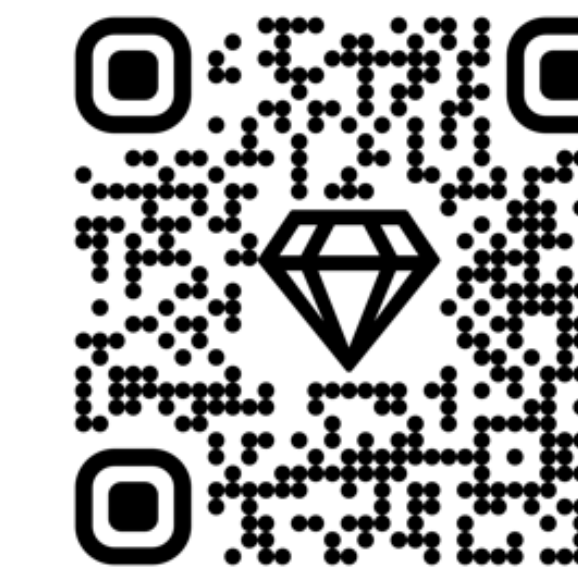
PERFORMANCE OF DEEP REINFORCEMENT LEARNING ALGORITHMS IN TWO-ECHELON INVENTORY CONTROL SYSTEMS

Francesco Stranieri^{†‡}, Fabio Stella[†] and Chaaben Kouki[§]

[†]Department of Informatics, Systems, and Communication, University of Milan-Bicocca, Italy

[‡]Department of Control and Computer Engineering, Polytechnic University of Turin, Italy

[§]Department of Operations Management and Decision Science, ESSCA School of Management, France



Introduction

This study addresses a *supply chain network optimization problem*, focusing on a two-echelon inventory system subject to capacity constraints. The goal is to determine the optimal quantity of products to produce and ship, minimizing total costs. Obtaining an optimal solution is challenging due to the *curse of dimensionality*.

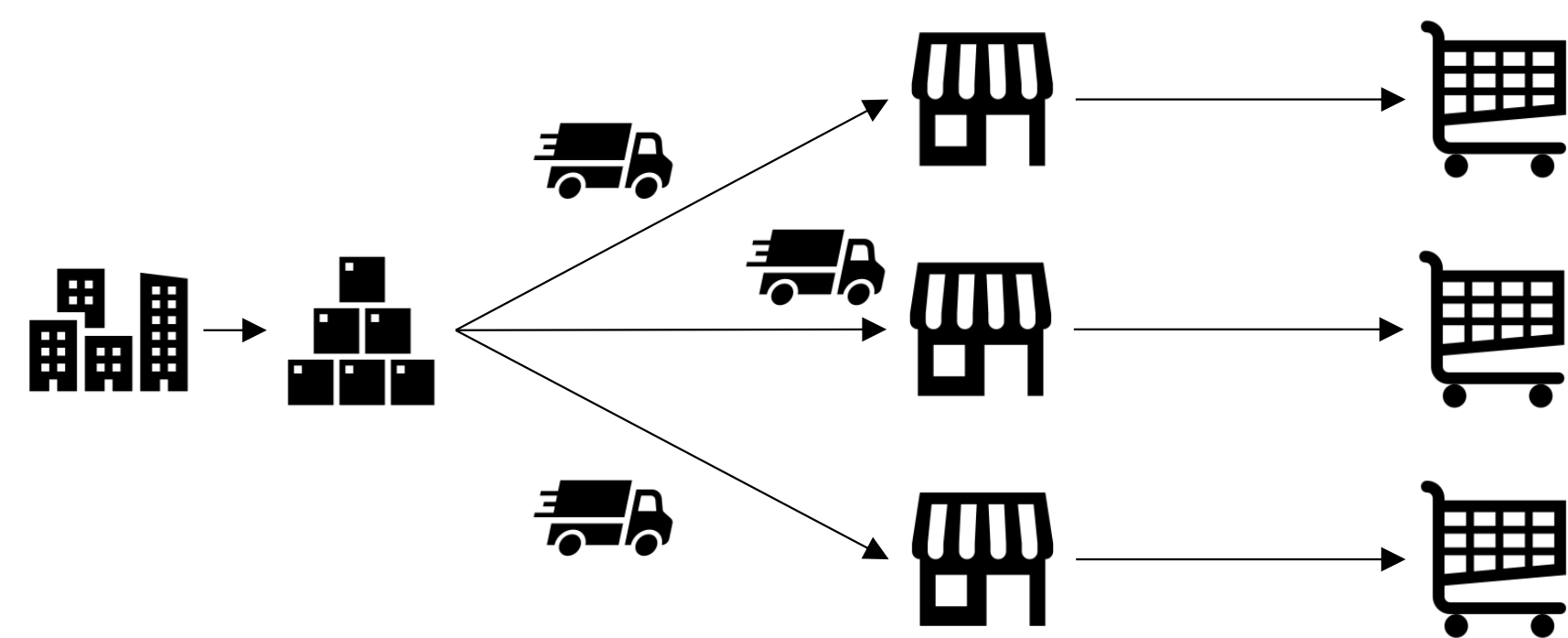


Fig. 1: A two-echelon inventory system composed of a central warehouse (first echelon) and three local warehouses (second echelon). Shopping carts represent customers' demands.

This study aims to evaluate the effectiveness of Deep Reinforcement Learning (DRL) algorithms by comparing their performance against traditional static inventory policies. This comparison will highlight the relative strengths of each approach, providing valuable insights to *decision-makers*.

Inventory Control System

For each product type $i \in \{1, \dots, I\}$, the factory determines its production level $a_{0,t}^i$, with a unit production cost of p_0^i . At each time step $t \in \{0, \dots, T\}$, $a_{j,t}^i$ units are shipped to local warehouse $j \in \{1, \dots, J\}$, which faces a stochastic demand $d_{j,t}^i$. Products shipped to local warehouse j are received after a lead time L_j , incurring a unit transportation cost of z_j^i . Each warehouse has a maximum capacity of c_j^i ($\sum_{i=1}^I c_0^i = c_0$), a storage cost of h_j^i per unit, and a stock level at time t of $q_{j,t}^i$.

Parameter	Explanation	Parameter	Explanation
I	Number of Product Types	p_0^i	Production Cost (per unit)
J	Number of Local Warehouses	z_j^i	Transportation Cost (per unit)
T	Episode Duration (time step)	$q_{j,t}^i$	Stock Level (units)
L_j	Transportation Lead Times (time step)	c_j^i	Storage Capacity (units)
$d_{j,t}^i$	Demand (units)	h_j^i	Storage Cost (per unit)
$a_{j,t}^i$	Production and Shipping Level (units)	b_j^i	Backorder Cost (per unit)

Tab. 1: The system parameters with their corresponding explanations (and units of measure). All these parameters are integrated and customizable into our open-source library, available on <https://github.com/frenkowski/SCIMAI-Gym>.

Products are non-perishable and supplied in discrete quantities. Any unsatisfied orders are back-ordered. Our study differentiates itself from existing literature [1] by incorporating *positive lead times* and considering *multiple product types*.

Main Contributions

We implement a continuous *action space* (i.e., the neural network directly generates the action value), denoted by \mathbf{a}_t , with an independent upper bound for each value, $0 \leq \sum_{i=1}^I a_{j,t}^i \leq c_j$. Additionally, our proposed balanced *allocation rule* ensures a fair distribution of products among local warehouses while maintaining non-negative stock at the central warehouse.

Algorithm 1 Our proposed allocation rule for continuous action spaces.

```

for all  $i \in \{1, \dots, I\}$  do
  if  $\sum_{j=1}^J a_{j,t}^i > q_{j,t}^i$  then
    while  $\sum_{j=1}^J a_{j,t}^i > q_{j,t}^i$  do
       $x \leftarrow \mathcal{U}(1, J)$ 
      if  $a_{x,t}^i > 0$  then
         $a_{x,t}^i \leftarrow a_{x,t}^i - 1$ 
    
```

The *state vector* includes all current stocks (on-hand and ordered but not yet delivered) and the last τ demand values (excluding d_t) and is denoted as:

$$(\mathbf{s}_t, \mathbf{d}_{1,t}, \mathbf{d}_{1,t-1}, \dots, \mathbf{d}_{J,t}).$$

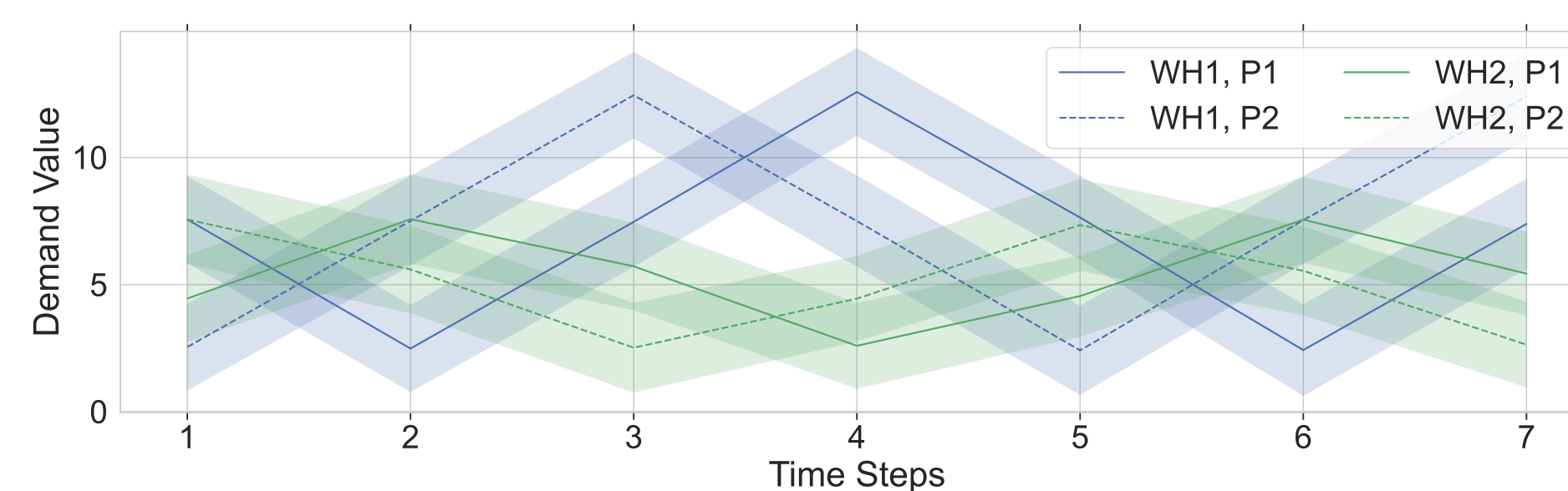


Fig. 2: An instance of the demand behavior considering two product types and two local warehouses.

Finally, we design the *reward function* for time step t as:

$$(R) \max \sum_{t=1}^T \sum_{j=1}^J \sum_{i=1}^I -[h_j^i(q_{j,t-1}^i + x_{j,t-L_j}^i - d_{j,t}^i)^+ + b_j^i(d_{j,t}^i - q_{j,t-1}^i - x_{j,t-L_j}^i)^+ + z_j^i a_{j,t}^i] - \sum_{t=1}^T \sum_{i=1}^I [h_0^i q_{0,t}^i + p_0^i a_{0,t}^i]$$

The first term captures the inventory cost; the second term represents the back-order cost, while the third denotes the transportation cost; the final two terms account for the inventory and production costs at the central warehouse.

References

- [1] Yimo Yan et al. "Reinforcement learning for logistics and supply chain management: Methodologies, state of the art, and future opportunities". In: *Transportation Research Part E: Logistics and Transportation Review* 162 (2022), p. 102712.
- [2] Francesco Stranieri and Fabio Stella. "A deep reinforcement learning approach to supply chain inventory management". In: *arXiv: 2204.09603* (2022).

Results

We conducted a set of numerical experiments to benchmark the performances of state-of-the-art *DRL algorithms* (i.e., PPO, PG, and A3C) against *static policies* (i.e., base-stock and (s, Q)-policies) [2]. We also implemented an *oracle* with perfect information (PI) on demand and an agent that knows the expected value of perfect information (EVPI) about demand. We performed 250 different simulations for each of the twelve experiments (by varying $I = \{1, 2\}$, $J = \{1, 2, 3\}$, and $L = \{1, 3\}$), evaluating seven demand realizations and reporting the *PI gap* achieved.

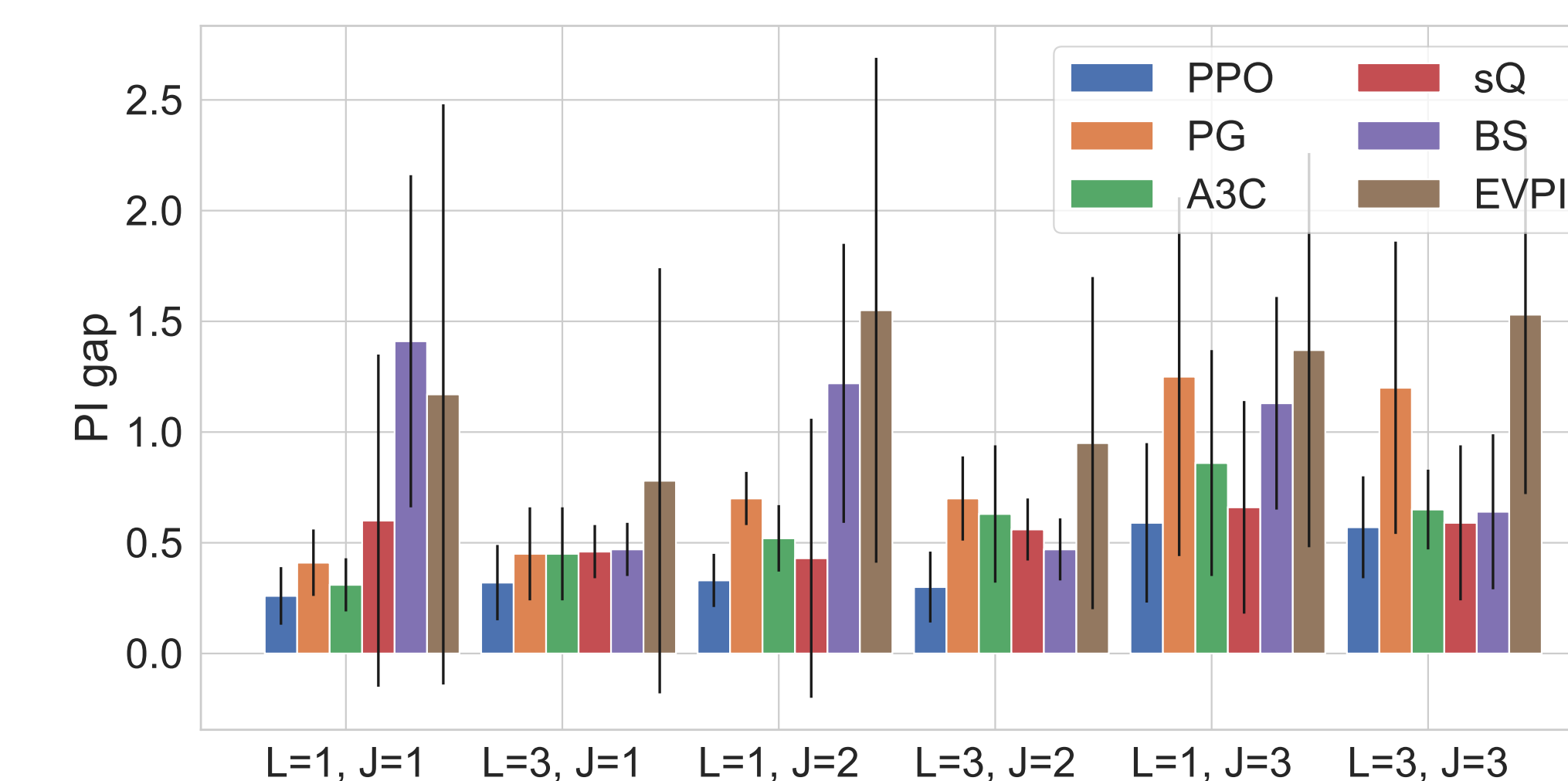


Fig. 3: With a single product type and lead times set to one, A3C outperforms PG, and the (s, Q)-policy demonstrates superior performance compared to the base-stock policy. When extending lead times to three, a degree of equilibrium is observed.

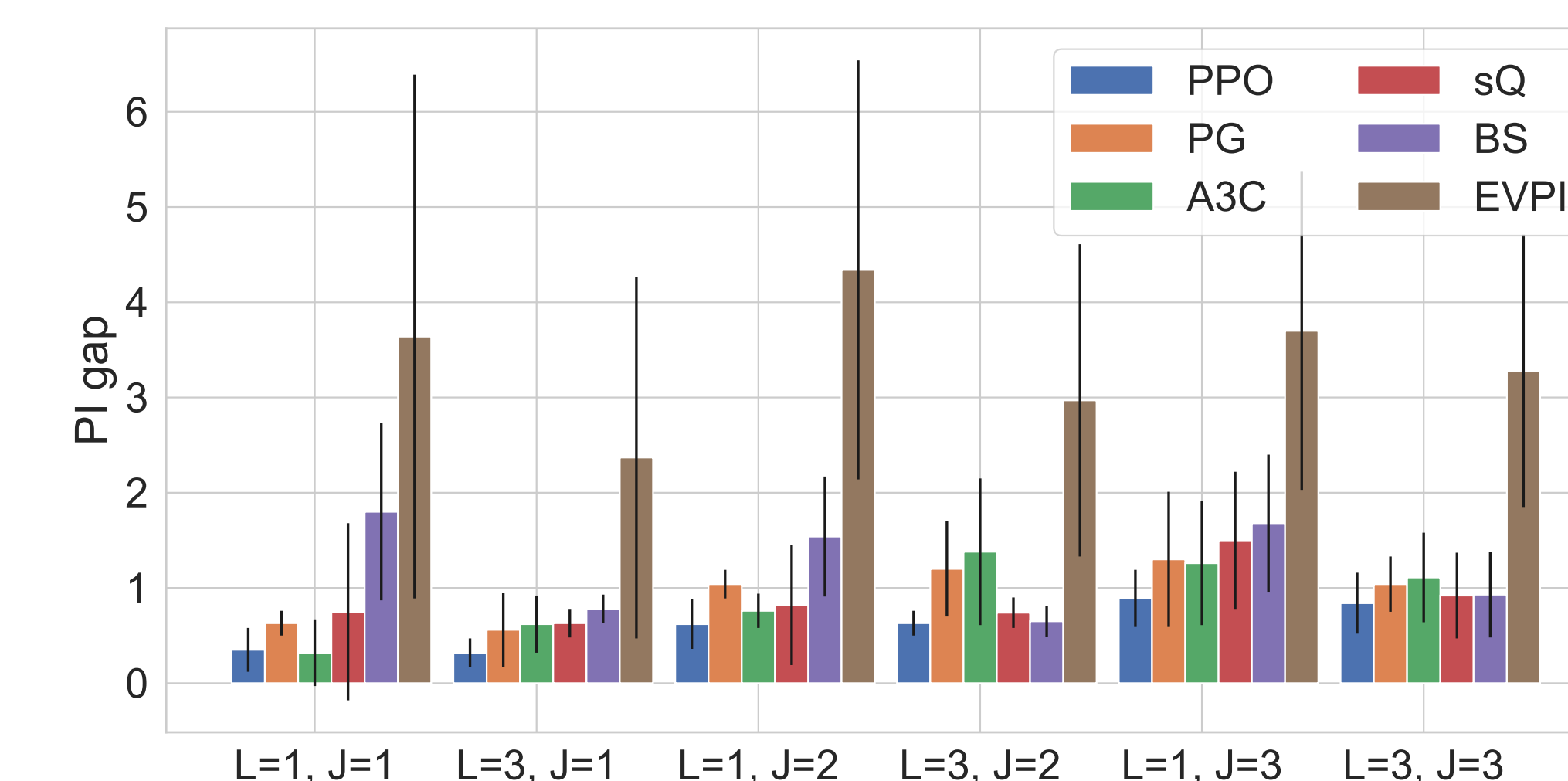


Fig. 4: With two product types and lead times set to one, DRL algorithms show slightly improved performance. When lead times are extended to three, static policies mildly outperform A3C and PG.

Our results indicate that *PPO consistently outperformed other algorithms* across all conducted experiments. The (s, Q)-policy performs commendably, emerging as the second-best choice. DRL algorithms also outperform EVPI, suggesting that complexities within the problem cannot be adequately addressed simply by considering the average demand value. Decision-makers should carefully consider these factors when selecting an *appropriate inventory management approach*.