

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/329718922>

Can Deep Reinforcement Learning Improve Inventory Management? Performance on Dual Sourcing, Lost Sales and Multi-Echelon Problems

Article in SSRN Electronic Journal · July 2019

DOI: 10.2139/ssrn.3302881

CITATIONS

12

READS

4,408

4 authors:



Joren Gijsbrechts

KU Leuven

5 PUBLICATIONS 19 CITATIONS

[SEE PROFILE](#)



Robert Boute

KU Leuven

65 PUBLICATIONS 570 CITATIONS

[SEE PROFILE](#)



Dennis J. Zhang

Washington University in St. Louis

32 PUBLICATIONS 284 CITATIONS

[SEE PROFILE](#)



Jan Albert Van Mieghem

Northwestern University

122 PUBLICATIONS 3,368 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Learning by Doing versus Learning by Viewing: An Empirical Study of Data Analyst Productivity on a Collaborative Platform at eBay [View project](#)



Forecasting Product Life Cycle Curves: Practical Approach and Empirical Analysis [View project](#)

Can Deep Reinforcement Learning Improve Inventory Management? Performance on Dual Sourcing, Lost Sales and Multi-Echelon Problems

Joren Gijsbrechts

Faculty of Economics and Business, KU Leuven, Belgium

Robert N. Boute

Vlerick Business School and KU Leuven, Belgium

Jan A. Van Mieghem

Kellogg School of Management, Northwestern University, United States

Dennis J. Zhang

Olin Business School, Washington University in St. Louis, United States

December 17, 2018; Revised July 29, 2019

The popularity of reinforcement learning is growing but is it effective in operations? We provide proof of concept that deep reinforcement learning (DRL) can be applied as a general-purpose technology to three classic, yet intractable inventory replenishment problems. Step-by-step guidance on how to apply DRL is proffered together with the code and a careful discussion of its performance, strengths and weaknesses.

Key words: artificial intelligence, deep reinforcement learning, inventory control, dual sourcing, lost sales, multi-echelon

1. Introduction

In the nexus of business and technology, no topic is hotter today than machine learning and artificial intelligence. We focus on deep reinforcement learning (DRL), the subfield of machine learning that develops policies for sequential decision-making problems. DRL employs deep neural nets to approximate the value or policy functions of Markov Decision Processes. This allows DRL to circumvent the curse of dimensionality, inherent to dynamic programming. With the help of DRL, programs have learned to play Atari games directly from image pixels [Mnih et al. 2015] and recently AlphaGo has beat the best human players in Go, the most complex board game in human history [Silver et al. 2017].

Despite the on-going frenzy about these breakthroughs, applications of DRL in industry contexts such as inventory management remain rather scarce. We demonstrate how and when DRL can be applied to classic, yet intractable inventory problems. The true strength of these general learning algorithms is that they provide a way to solve a diversity of problems rather than relying on extensive domain knowledge or restrictive assumptions. In the inventory management literature a variety of model-specific heuristics exist; it typically depends on the model assumptions which heuristic performs best. As a “general purpose technology”, DRL is a single, recognizable generic technology that can serve many different uses. We establish proof-of-concept with rigorous performance benchmarks to dual-sourcing or dual-mode, lost sales and multi-echelon inventory models.

Under dual sourcing, inventory can be replenished from a fast but expensive source and from a regular, cheaper source with longer lead time. In the equivalent dual-mode problem, inventory can be replenished from a single source using two complementary transportation modes. In lost sales inventory models customers walk away or receive an expedited shipment at a premium cost when the desired product is not available (and hence no backlogging is allowed). Multi-echelon inventory management includes additional stages or echelons of the supply chain that can hold inventory, such as an intermediate central warehouse between an upstream plant and downstream retailers.

These conceptually-simple problems have vexed supply chain management scientists for decades. Little is known about the optimal policy structure and traditional solution methodologies such as dynamic programming quickly become intractable due to the curse of dimensionality: the state-space grows exponentially as replenishment lead times get longer. Heuristic policies are typically problem-dependent and rely on restrictive assumptions, limiting their use in different settings. We explore whether a generic technology like DRL provides an alternative.

We aspire to make the following contributions: (1) Provide proof-of-concept that deep reinforcement learning can tackle a variety of inventory problems that have intractable dynamic programs; (2) Numerically evaluate performance relative to the optimal dynamic program and compare the performance and ease-of-use versus heuristic policies and approximate dynamic programming methods; and (3) Provide a tutorial on how, why and when DRL works.

2. Selected Literature reviews

We provide a brief introduction on well-performing heuristics on dual sourcing (Section 2.1), lost sales (Section 2.2) and multi-echelon (Section 2.3) inventory models. These heuristics have been developed and shown to perform well under the conventional assumptions: i.i.d. demand, linear sourcing costs and deterministic lead times. As soon as one of these assumptions is relaxed or other practice-specific constraints are added, the performance of these heuristics may suffer and adjusted or new well-performing policies are desired. This is precisely where approximation methods provide a solution, which we review in Section 2.4.

2.1. Selected Literature Review - Dual Sourcing Inventory Models

Dual sourcing or dual-mode replenishment has been studied since the pioneering inventory management models of the mid-twentieth century. Fukuda [1964] proved that the optimal policy follows a base-stock structure when the lead time difference between both sources is exactly one period. When the lead time difference exceeds one period, however, Whittemore and Saunders [1977] showed that no simple structure prevails and the optimal policy depends on the vector of outstanding orders. As costs related to replenishments within the expedite lead time do not impact the optimal policy [Sheopuri et al. 2010], the optimal policy thus depends on the $(l^r - l^e)$ outstanding

receipts, in which l_r and l_e respectively represent the lead time of the regular slow and expedite supplier. This inspired the development of various heuristic policies that work around this complexity by collapsing the state vector to one or two inventory positions. For instance, the dual-base stock policies of Scheller-Wolf et al. [2003] and Veeraraghavan and Scheller-Wolf [2008] aggregate the pipeline inventory: the *Single-Index (SI)* policy uses one inventory-position over the long lead time while the *Dual-Index (DI)* uses both the long and short lead time inventory positions. The *Capped Dual Index (CDI)* policy [Sun and Van Mieghem 2018] extends the dual index policy by introducing a cap on the order to the slow source. The *Tailored Base-Surge (TBS)* policy [Allon and Van Mieghem 2010] places a constant order to the slow source, and operates with a base-stock policy to determine orders at the fast source. Xin and Goldberg [2017] prove that the TBS policy is asymptotically optimal when the lead time difference between both sources grows to infinity.

Sheopuri et al. [2010] show the non-optimality of order-up-to-policies for non-consecutive lead times and propose policies making better use of the structure of the pipeline inventory vector: the *Weighted Dual Index* policy weighs different elements of the inventory vector prior to ordering up to the base-stock levels while *Vector Base-Stock (VBS)* policies employ a base-stock for the expedited source and a vector base-stock policy identical to Zipkin [2008a] for the slow source. Hua et al. [2015] show that the dual sourcing value function satisfies $L^\frac{1}{2}$ -convexity (read *el-natural*), similar to the lost sales inventory model [Zipkin 2008b]. Inspired by this property, they show that the optimal orders to the slow source are more sensitive to the longest outstanding orders, and the optimal orders to the fast source depend more on the soon-to-arrive outstanding orders. By developing upper and lower bounds on the slow source's order and using a weighted average of these bounds to determine the slow source's orders, they generalize the VBS policy to the *Best Weighted Bounds* policy. Chen and Yang [2019] exploit the $L^\frac{1}{2}$ -convexity property to develop a linear programming approximate dynamic programming (LP-ADP) algorithm to dual sourcing with supply uncertainty. We elaborate further on LP-ADP in Section 2.4.

2.2. Selected Literature Review - Lost Sales Inventory Models

Arrow and Karlin [1958] show that in a single sourcing backlogging model with constant lead time, inventory on hand and outstanding orders can be collapsed into one single number (i.e. the inventory position) and that a single-dimensional base-stock policy is optimal. They also indicate that this is no longer the case when demand is lost instead of backlogged. Instead, the optimal policy becomes dependent on the full inventory pipeline vector for non-zero lead times. Similar to dual sourcing, the complexity then grows exponentially in the lead time. More structural properties and bounds are provided in Morton [1969]. Sheopuri et al. [2010] show the similarity between the lost sales and the dual sourcing models. This is also apparent in the type of heuristic policies for

both models. The *Constant Ordering* policy [Reiman 2004], for instance, is a Tailored Base-Surge policy with only constant replenishment from a single source, that is asymptotically optimal for long lead times in lost sales inventory models [Goldberg et al. 2016].

Base-Stock policies generally perform poorly in lost sales inventory models [Zipkin 2008a], except for large penalty costs, in which case optimal inventory levels are high and stock-outs are rare. Huh et al. [2009] prove that base-stock policies are asymptotically optimal for large penalty costs. Xin [2019] combine the two asymptotic results (i.e. constant ordering for increasing lead times and base-stock policies for increasing penalty costs) in the *Capped Base-Stock* policy, inspired by its Capped Dual Index counterpart in dual sourcing [Sun and Van Mieghem 2018].

The *Myopic* policy generally performs well for the lost sales inventory model [Morton 1971, Zipkin 2008a]: it computes in each period t the order quantity that minimizes the expected costs in the period when the order arrives (i.e. in period $t + l$, with l representing the lead time) by recursively computing the expected stock levels in period $t + l$. These findings are more rigorously confirmed by Brown and Smith [2014] who use the value of perfect information in combination with the $L^\frac{1}{2}$ -convexity property (which was proven by Zipkin [2008b]). They develop tight lower bounds and show how myopic policies consistently perform close to optimality on the lost sales problem even in settings where dynamic programming is intractable. Extensions such as the 2-period myopic policy [Zipkin 2008a] that considers two periods perform even better. Further increasing the myopic horizon comes at the expense of increased computation time (similar to the dynamic program in which the infinite horizon is used to choose actions). The one- and two-period myopic policies are essentially approximate dynamic programming methods as they require, at each time epoch, the iteration (or evaluation) over all states that can be reached during the considered horizon [Xin 2019]. They thus also suffer from the curse of dimensionality and are more computationally intensive than base-stock or constant order policies.

Inspired by Zipkin's (2008b) finding that the value function of the lost sales problem is $L^\frac{1}{2}$ -convex, Sun et al. [2016] develop an LP-ADP algorithm, identical to that of Chen and Yang [2019] for dual sourcing with supply uncertainty.

2.3. Selected Literature Review - Multi-Echelon Inventory Models

Multi-echelon models include multiple stages or echelons in the supply chain that can hold inventory. For instance, central warehouses can pool inventory prior to allocating it to downstream retailers. Seminal works such as Clark and Scarf [1960] and Federgruen and Zipkin [1984] characterize the optimal policy structure under several strong assumptions such as, respectively, having a serial system (i.e. no parallel stocking facilities) or no option to stock at intermediate terminals. de Kok et al. [2018] provide an extensive review of the variety of multi-echelon inventory models studied based on number of echelons, network structure, holding cost functions, etc.

We study the multi-echelon inventory model employed in Van Roy et al. [1997]. They adopt a neuro-dynamic programming approach to solve a two-echelon model with one warehouse and multiple retailers. The model is a hybrid between backlogging and lost sales, as described in Section 7. Closely related to this model are the partial lost sales and divergent network models studied in Nahmias and Smith [1993, 1994].

2.4. Selected Literature review - Approximate Dynamic Programming

Markov Decision Processes (MDPs) provide a mathematical framework to solve sequential decision-making problems with a countable state and control space [Puterman 1994]. Traditional solution methodologies such as linear programming (LP) or dynamic programming (DP) using either value or policy iteration, generally do not scale well to large problem sizes. Therefore, several *approximate dynamic programming* (ADP) methods have been developed as described in the seminal book by Powell [2011]. We provide a selective review, starting from methods that heavily exploit the problem structure to the most general purpose technologies available to date.

One branch of ADP methods exploits the problem structure. For example, for inventory models with zero or one period lead time, Halman et al. [2009] and Chen et al. [2014] avoid the need to iterate over the entire state space, speeding up the DP optimization and allowing the algorithm to solve close to polynomial time. Because of their heavy reliance on the problem structure, Chen et al. [2014] discuss the desire of more general purpose models. Our work may address this desire.

A second branch of ADP methods reduces the size of the problem by aggregating states, e.g. by clustering states based on features. Fang et al. [2013] and Giannoccaro and Pontrandolfo [2002] cluster the pipeline inventory into subsets by aggregating parts of the pipeline inventory. Based on their expertise into the specific problem, Van Roy et al. [1997] manually select 23 features from a multi-echelon supply chain to cluster the state space. These features include aggregations of parts of the inventory vectors, the variance among retailer inventory levels, products of inventory levels etc. It is often left up to the discretion of the user how to choose good cluster sizes, which can be a formidable challenge. Keller et al. [2006] automate this state aggregation process by using a neighborhood component analysis. While state aggregation can speed up computation, important information can get lost. This is always the case when the optimal policy depends on each element of the state, as in the three inventory problems that we study.

A third branch approximates the value or policy functions of MDPs to generate near-optimal policies. Selecting a good approximation is challenging if little is known about the structure of the optimal value or policy function. In inventory management, various function approximations are used, such as a linear combination of the state variables [Keller et al. 2006] or fitting a specific quadratic function on the state variables [as in Chen and Yang 2019, Sun et al. 2016]. The latter

exploits the L_1 -convexity property of the lost sales and dual sourcing model with supply uncertainty to fit a specific L_1 -convex quadratic shape, as introduced by Murota [2003]. As it is often infeasible to iterate or explore the entire state space when approximating the value function, efficient sampling of states is required. Some methods rely on well-known heuristics to sample states: Chen and Yang [2019] leverage the Single/Dual Index and Tailored Base-Surge policies, and Sun et al. [2016] use the myopic lost sales policy to sample states. Both then apply linear programming to fit the approximation coefficients of the value functions using the sampled states.

Linear programming approximate dynamic programming (LP-ADP) falls in the third branch of ADP techniques. It assumes a specific functional approximation of the sample states using a well-known heuristic and uses linear programming to obtain the coefficients of these value functions. This approach circumvents the computational explosion characterizing traditional linear programming to solve MDPs: the need to include the value function of each state in the objective while each state-action pair requires a unique constraint. Instead, only the sampled states are included in the objective function and only the sampled state-action pairs are included in the constraints. Solving the LP then provides the coefficients of the approximating value functions, which then provides a vehicle to evaluate the value function at unvisited states. An “LP-greedy” policy minimizes the current cost plus the cost-to-go of the next stage. Both Chen and Yang [2019] and Sun et al. [2016] report small savings in comparison to the sampling policies. We employ LP-ADP as a benchmark in our numerical experiments.

Reinforcement learning (RL) is a mathematical framework to solve MDPs without requiring an exact representation of the environment. A RL agent *learns* to maximize expected rewards by interacting with an environment. Mathematically, RL algorithms develop a good approximation of the value or policy function of the underlying MDP. Choosing a good functional approximation and efficient sampling of states and actions is key. While there are RL algorithms that heavily exploit problem structure, we focus on the recently developed general purpose RL algorithms. A major benefit of these RL algorithms in comparison to problem-specific heuristics is that they can *learn* policies in environments with less stylized or non-conventional assumptions. This is valuable in settings where off-the-shelf inventory replenishment heuristics are not available or do not perform well. A well-performing generic technology then adds value compared to tailor-made policies for one specific company or business.

Although classic RL algorithms, such as temporal difference and Q-learning, had some success in the past [Tesauro 1995, Kohl and Stone 2004, Watkins 1989], original RL approaches still suffered from the lack of scalability [Strehl et al. 2006]. Therefore, in recent years, reinforcement learning is combined with deep learning (i.e., powerful function approximation using deep neural networks). Deep reinforcement learning (DRL) thus uses deep neural networks to approximate the value or

policy functions of MDPs. Although both RL and DRL have been studied extensively in the Computer Science and Operations Research literature [Sutton and Barto 1998, Mnih et al. 2015], the accessibility of more computational power and recent algorithmic breakthroughs have sparked a new interest in DRL. The Asynchronous Advantage Actor-Critic (A3C) algorithm that we will employ is one of the most popular recent DRL algorithms due to its excellent performance and fast training time [Mnih et al. 2016]. Moreover, we show that with only minor modifications, it is capable to develop good policies on three intractable inventory problems.

To the best of our knowledge, only Oroojlooyjadid et al. [2017] and Oroojlooyjadid et al. [2016] have applied the latest deep reinforcement learning models to inventory management problems, respectively to the Beer Distribution Game and the newsvendor problem. While our paper also applies DRL to inventory problems, it differs from those working papers in two distinctive ways. First, we provide optimality gaps on three different inventory problems whose stochastic dynamic program become quickly intractable. Second, we show the versatility of DRL by applying it to three different inventory problems with limited modification of the algorithm, thereby demonstrating that DRL resembles a general purpose technology.

3. DRL Solution Approach to Dual Sourcing Inventory Replenishment

Deep Reinforcement Learning provides a way to approximate and solve large Markov Decision Processes for which traditional dynamic programming methods are intractable. We will demonstrate in this section how the Asynchronous Advantage Actor-Critic (A3C) algorithm operates. We use the dual sourcing model as the running example and show in later sections how both the lost sales and multi-echelon can be plugged into this framework, showing the versatility of the algorithm.

Consider the periodic-review inventory replenishment of a single item over an infinite horizon. In the conventional dual sourcing or dual mode model, inventory can be replenished at unit cost c_r using a regular source with lead time l_r and using an expedited source with lead time $l_e < l_r$ at premium unit cost $c_e > c_r$. At the beginning of any period t , two order quantities, $q_t^r \geq 0$ and $q_t^e \geq 0$, must be decided knowing the last observed inventory on hand, I_{t-1} , and outstanding receipts or “pipeline” vector $Q_{t-1} = (q_{t-l_r}^r + q_{t-l_e}^e, q_{t-l_r+1}^r + q_{t-l_e+1}^e, \dots, q_{t-l_r+l_e-1}^r + q_{t-1}^e, q_{t-l_r+l_e-2}^r, \dots, q_{t-1}^r)$ ¹. Then, orders $q_{t-l_r}^r$ and $q_{t-l_e}^e$ are received and added to the on-hand inventory. Finally, the unknown demand d_t is realized and subtracted from the on-hand inventory. Excess demand is backlogged so that the inventory and pipeline vector evolve as: $I_t = I_{t-1} + q_{t-l_r}^r + q_{t-l_e}^e - d_t$ and $Q_t = (q_{t-l_r+1}^r + q_{t-l_e+1}^e, q_{t-l_r+2}^r + q_{t-l_e+2}^e, \dots, q_{t-l_r+l_e}^r + q_t^e, q_{t-l_r+l_e-1}^r, \dots, q_{t-1}^r)$.

The dual sourcing problem can be modeled as a Markov Decision Process with state $s_t = (I_{t-1}, Q_{t-1})$ at time t . Let \mathcal{S}_t denote the set of admissible states at time t . After taking action

¹This formulation assumes strictly positive lead times. If $l_e=0$, there are no outstanding expedited orders.

$a_t = (q_t^r, q_t^e)$, the cost c_t incurred in period t consists of sourcing costs, per-period holding cost h for each unit held in inventory-on-hand or per-period backlog cost b for each unit backlogged:

$$c_t(s_t, a_t) = q_t^r c_r + q_t^e c_e + h[I_t]^+ + b[I_t]^- . \quad (1)$$

Let \mathcal{A}_t denote the set of admissible order policies $\pi_t = \{a_{t+j} : j = 0, 1, \dots\}$. The objective is to find an admissible order policy that minimizes the expected present value \mathcal{C}_t of future costs. Assuming henceforth sufficient regularity including countable, compact state and action spaces, this cost when starting from state s_t and following policy π_t is:

$$\mathcal{C}_t^{\pi_t}(s_t) = \sum_{j=0}^{\infty} \gamma^j \mathbb{E}^{\pi_t} c_{t+j}(s_{t+j}, a_{t+j})$$

where $\gamma \in (0, 1)$ is the discount factor and \mathbb{E}^{π_t} is the expectation operator when following policy π_t . Define the optimal value function $v(s_t)$ as the infimum of $\mathcal{C}_t^{\pi_t}(s_t)$. Assuming sufficient regularity and countable states and actions, the value functions are obtained when following an optimal policy and solve the celebrated recursive *Bellman equations* [Bellman 1954]:

$$v(s_t) = \min_{a_t \in \mathcal{A}_t} \left\{ c_t(s_t, a_t) + \gamma \sum_{s' \in \mathcal{S}_{t+1}} \mathbb{P}(s_{t+1} = s' | s_t, a_t) v_{t+1}(s') \right\} .$$

Despite their simple appearance, solving these Bellman equations can quickly become problematic as their problem complexity suffers from the triple curse of dimensionality [Powell 2011], which is driven by the dimension of (1) the state space $(l_r - l_e)^2$; (2) the action space (all possible sourcing quantity combinations a_t); and (3) the transition probability matrix (dependent on the size and number of possible demand realizations).

To cope with the curse of dimensionality, the dimensions of the problem can be reduced by applying approximate dynamic programming. As discussed in the literature review, one ADP branch includes *reinforcement learning*. In their simplest form, value-based methods such as tabular Q -learning [Watkins 1989] store estimated values of $v(s)$ in lookup tables. In Q -learning, *action-value functions* $Q(s_t, a)$ represent the future expected cost of taking action a from state s_t and then following the best-known policy from state s_{t+1} . Note that the Q -values only differ in the first step from the value function $v(s)$. These Q -values can iteratively be improved: $Q_{t+1}(s_t, a_t) = (1 - \alpha)Q_t(s_t, a_t) + \alpha(c_t + \gamma \max_{a \in \mathcal{A}} \mathbb{E}Q(s_{t+1}, a))$, in which α represents the learning rate that trades off how much the new observation $c_t + \gamma \max_{a \in \mathcal{A}} \mathbb{E}Q(s_{t+1}, a)$ is used to update the initial estimate $Q_t(s_t, a_t)$. New states and actions can be explored by following ϵ -greedy methods that trade off following the best-known policy, with probability $(1 - \epsilon)$ against exploring new actions, with probability ϵ .

² While (I_{t-1}, Q_{t-1}) has dimension $l_r + 1$, future costs in periods t to $t + l_e - 1$ cannot be influenced and are sunk at time t . Hence the state space dimension of the dynamic program can be reduced to $l_r - l_e$ [Sheopuri et al. 2010].

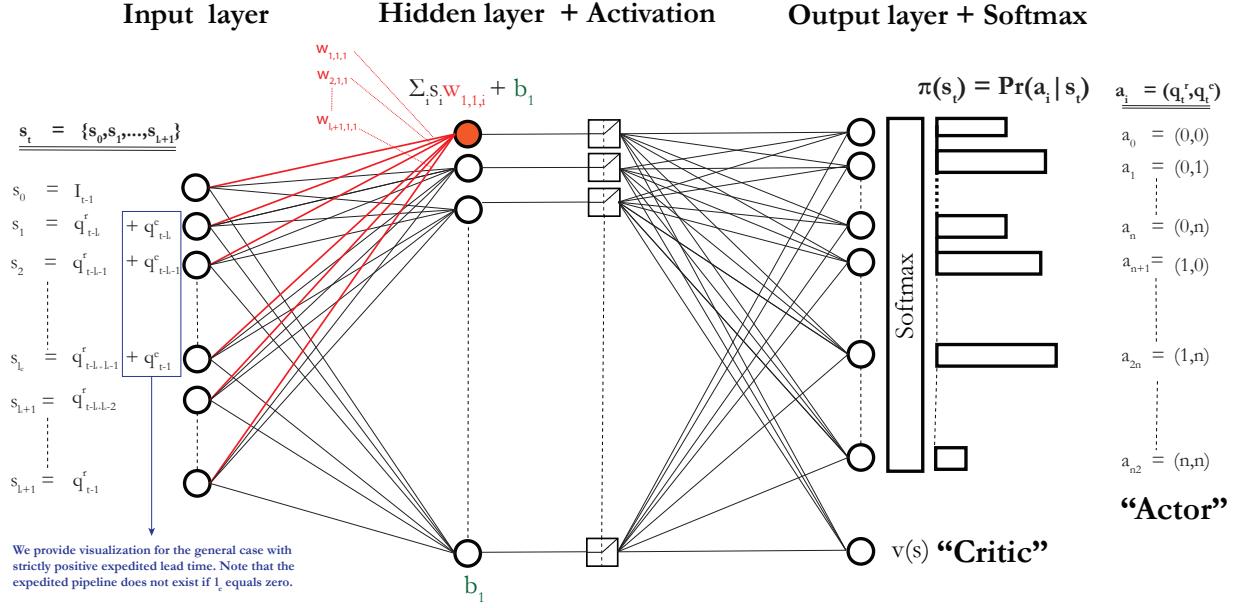


Figure 1 **Visualization of a simple neural net of the A3C where the input dimension is $l_r + 1$, one hidden layer is used and n^2 actions can be chosen. The output consists of (1) a value $v(s_t)$ (estimated by the “critic”) when using (2) a stochastic policy over the action space $\pi(s_t) = \mathbb{P}(a_t | s_t)$ (estimated by the actor using a softmax function on the last layer that normalizes all values into a probability vector).**

More sophisticated value-based methods use parametric approximations of the values $v(s)$ to circumvent the curse of dimensionality when storing all function values in lookup tables. Neural networks are an example of approximators of non-linear functions. Neural nets have an input layer, several hidden and activation layers in sequence, and an output layer. Figure 1 illustrates a neural net with one hidden layer. Its input is the $(l_r + 1)$ -dimensional state vector $s_t = (s_0 = I_{t-1}, s_1 = q_{t-l_e}^r + q_{t-l_e}^e, s_2 = q_{t-l_e+1}^r + q_{t-l_e+1}^e, \dots, s_{l_r} = q_{t-l_e+l_r-1}^r + q_{t-l_e+l_r-1}^e, s_{l_r+1} = q_{t-1}^r)$. A layer’s output is a linear function of its input followed by a non-linear “activation layer.” The value of the first node of the hidden layer is $\sum_{i=0}^{l_r+1} w_{1,1,i} s_i + b_1$ where the weights of the red edges ($w_{1,1,1} \dots w_{1,1,l_r+1}$) and the bias b_1 are tuning parameters. The activation layer adds non-linearity to the linear function. For example, a positive-part operator (called “rectified linear unit” ReLu) only passes positive node values to the next layer. Concatenating layers provides powerful approximations to non-linear value functions. Traditional neural networks use mostly one or two layers while recent *deep learning* employs a larger number of hidden layers.

In contrast to value-based methods, *policy-based* methods directly develop a policy. A deterministic policy prescribes for each state a single action. A stochastic policy prescribes for each state a probability distribution over all actions that specifies the probability that each action should be taken.

At the juncture of value-based and policy-based methods, actor-critic methods combine both techniques by relying on an *actor* who develops a policy that is evaluated by a *critic*. We will explore the power of one such actor-critic method, namely the Asynchronous Advantage Actor-Critic (A3C) algorithm. In what follows, we use notation θ to define all parameters that can be updated during training of the model (i.e. all weights of the edges and biases of each layer). Given set θ , A3C generates for every state s_t a stochastic policy: $\pi(s_t; \theta) = \mathbb{P}(a_t | s_t)$ (i.e. a probability distribution over all actions) and one value function $v^\pi(s_t; \theta)$ (representing the expected future discounted cost of following the stochastic policy). In our experiments, we use a fully connected network, implying that all nodes between layers are connected and the actor and critic use the same hidden layers. This delivered good performance on our dual sourcing problem, but it does not necessarily need to be the case. The output of the actor passes a “softmax” function that normalizes the output layer values into a proper probability vector that defines the stochastic policy, as shown in Figure 1.

Training the A3C algorithm involves adjusting the weights and biases θ of the neural net to improve the policy of the actor and the estimation of the value function of the critic. The A3C algorithm employs a special structure where n parallel learners together, yet asynchronously, update the parameters θ as follows: Each learner $i \in \{1, \dots, n\}$ interacts independently with its own copy of the environment and uses its own neural net with parameters θ^i . These learners are used to update the parameters of the ‘global’ neural net θ^g . The agents learn and update the global network through a training *buffer* of m periods of observed states, actions and costs as follows: During the k -th buffer, the current policy π_k is simulated for m subsequent periods starting at state $s_{m(k-1)}$. During the simulation, each parallel learner keeps track of the incurred costs $(c_{m(k-1)}, \dots, c_{mk})$, the actions taken $(a_{m(k-1)}, \dots, a_{mk})$, the visited states $(s_{m(k-1)}, \dots, s_{mk})$ and the resulting state at the end of this episode (s_{mk+1}) . Given these records, each learner computes its loss function for training iteration and buffer k , denoted by Loss_k and to be defined in detail soon. Each learner i then computes the gradients of the total loss with respect to the parameters θ^i of its own local neural net and iteratively adjusts the parameters of the global net θ^g using a stepsize $\alpha > 0$ in the direction of the gradient to reduce the loss:

$$\theta_{k+1}^g = \theta_k^g - \alpha \nabla_\theta^i \text{Loss}_k / \|\nabla_\theta^i \text{Loss}_k\|.$$

This process happens *asynchronously*, meaning that updates to the weights of the global network are made immediately whenever an agent simulated its buffer and computed its gradient. The weights of the global network are thus continuously updated by using the gradients from the local learners. Once the global network is updated, the weights of the global network are then copied back

to the local learner. The motivation for using parallel learners is that each learner is interacting with its own copy of the environment, which reduces the probability of over-fitting the global neural net on specific state regions. To avoid exploding gradients, a clipping function bounds the gradients. The gradient-descent step size depends on the used optimization method. We employ Adam, developed by Kingma and Ba [2015], which dynamically adapts the step size based on the improvement of past updates.

The total loss function used is the sum of three specific losses that we now define:

$$\text{Loss} = \text{Value Loss} + \text{Policy Loss} + \text{Entropy Loss}.$$

The value loss measures the quality of the value function approximation by the difference between the future discounted cost and the approximated value function at each observed state in the episode buffer. The future discounted cost of a state $s_p = s_{mt-k}$ in the t -th episode buffer consists of the costs observed k steps until the end of the episode buffer, $C_p^{(k)} = \sum_{i=0}^k \gamma^i c_{p+i}$, plus the infinite-horizon discounted costs after the episode buffer, which is approximated by the value function in the last period of the buffer, $\gamma^{k+1} v^{\pi_t}(s_{tm+1}; \theta)$. The value loss is then defined as the sum of the squared errors of all states $(s_{m(t-1)}, \dots, s_{mt})$ within the episode buffer. We use a value regularization term β_V (which we fixed at 0.25) to prevent over-fitting:

$$\text{Value Loss} = \beta_V \sum_{p=0}^m \left(-v^{\pi_t}(s_{m(t-1)+p}; \theta) + C_{m(t-1)+p}^{(m-p)} + \gamma^{m-p} v^{\pi_t}(s_{mt+1}; \theta) \right)^2.$$

The policy loss is used to improve the policy by selecting “better” actions. For each action taken at period p during the episode t ($a_p(s_p)$), we compute the difference between the cost of this action plus the discounted value of the next state ($c_t(s_t, a_t) + \gamma v^{\pi_t}(s_{t+1})$) and the value function at the current state ($v^{\pi_t}(s_t)$). This difference can be positive, which means that the action is not optimal using current value function, zero, or negative, which means the current action is better than the action suggested by the policy. Therefore, minimizing this difference helps to choose better actions with respect to the current value function. We use Generalized Advantage Estimation³ [Schulman et al. 2015], in which the observed difference within the episode buffer are once again discounted. Since, at a given state, each action is only taken with a certain probability, the loss is corrected by multiplying by the log of the probability of selecting action a_t :

$$\text{Policy Loss} = \sum_{i=0}^m \left(\log \mathbb{P}(a_{m(t-1)+i} | s_{m(t-1)+i}) \sum_{p=0}^k \gamma^p (c_{m(t-1)+i} + \gamma v^{\pi_t}(s_{m(t-1)+i+1}; \theta) - v^{\pi_t}(s_{m(t-1)+i}; \theta)) \right).$$

³ Generalized advantage estimation uses an additional weighing parameter $\lambda \in [0, 1]$ which we keep fixed at one as it provides good performance in our setting.

Finally, the entropy loss is used to avoid that the A3C algorithm over-fits and converges to local optima. Minimizing entropy ensures that the model keeps the right balance between exploring new actions and exploiting actions that are known to perform well. The entropy for each episode buffer is defined as the entropy of the probability function over actions taken, which is the logarithm of the probability mass function over actions taken: $\mathbb{P}(a_i|s_i) \log \mathbb{P}(a_i|s_i)$. This entropy is minimized when all actions have the same probability and is equivalent to encouraging the algorithm to explore new actions. An entropy regularization term β_E determines how much exploration is added to the loss function:

$$\text{Entropy Loss} = \beta_E \sum_{p=0}^m \mathbb{P}(a_{m(t-1)+p}|s_{m(t-1)+p}) \log \mathbb{P}(a_{m(t-1)+p}|s_{m(t-1)+p}).$$

4. How to tune a DRL algorithm?

Thanks to the open nature of the machine learning community, many deep reinforcement learning algorithms are freely available. While this should greatly facilitate their application, it is our experience that the difficult tuning process often demotivates researchers to explore the full potential of these new methods. Indeed, determining the suitable hyperparameters of the A3C algorithm (summarized in Table 1) is a non-trivial and very time-consuming task because evaluating hyperparameter sets is expensive and noisy. Therefore, we share our experience on how to achieve good performance for three different inventory problems, together with our code⁴. This may help other researchers to directly build on our experience when applying DRL to other operations problems.

For the dual sourcing model, we initially tested various hyperparameter settings and kept track of well-performing hyperparameter values. To further improve the results, we stored the best models (i.e., trained neural networks) and restarted from these models using different entropy factors and learning rates. This further improves already well-performing models and uses less exploration and more exploitation as training progresses. This labor intensive tuning process motivated us to develop an automated tuning strategy without manual intervention, which we also tested for the lost sales and multi-echelon inventory models.

Building on our experience of manually tuning various model instances of the dual sourcing problem, we developed the following automated tuning process for the A3C algorithm applied to inventory problems: First, we decided to keep the following hyperparameters fixed: four layers in the neural network with widths 150, 120, 80 and 20, respectively; each layer followed by a ReLu activation; value regularization 0.25; four parallel learners, and clipping rate 40. These hyperparameters did perform well and any manual tuning deviations did not result in any significant improvement.

⁴ A link to the code will be shared along with the publication of this paper.

Hyperparameter	Well-performing values	Range for automated tuning
Number of layers	4	
Width of layers	[150,120,80,20]	
Value regularization (β_V)	0.25	
Activation functions	ReLU	
Number of parallel learners	4	
Clipping rate	40	
Learning rate	10^{-4}	$[10^{-7}, 10^{-2}]$
Entropy regularization (β_E)	10^{-7}	$[10^{-10}, 10^{-2}]$
Buffer size	20	[20,100]

Table 1 Hyperparameters of the A3C algorithm that require tuning. The first column contains one set that performs well across all settings. The second column contains the ranges used in the automatic tuning process.

Equivalently, the computational effort to optimize those parameters did not outweigh their benefits in our experience. Therefore, they provide a good starting point. Of course, it is possible that changes to these hyperparameters can perform equally well or better in other settings.

Second, we took 200 random sample points of the remaining three hyperparameters—i.e. the length of the buffer, the entropy regularization rate and the learning rate—over pre-defined ranges. We determined those ranges by monitoring the three components of the loss function (value, policy and entropy loss) during manual training. This provided insight into the relevant ranges of the three parameters where the A3C algorithm tended to converge to good results. We evaluated the expected cost of the A3C policy by simulating 10 sample paths of 100,000 periods.

Third is the final choice of the three hyperparameter values. The simplest solution is to pick the sample point with the best results. We went one step further and fitted a quadratic function, specifically an 3D ellipsoid, through the lower convex hull of the 200 points to provide confidence in the existence of a global minimum of the tuning parameters. Later, we will visualize these ellipsoids.

Table 1 summarizes both the values of hyperparameters that are kept fixed as well as the ranges used in the automatic tuning of the remaining three hyperparameters. Note that the learning rate and entropy regularization are sampled from a continuous interval, while the buffer size is constrained to be integer. Once the hyperparameters are chosen, the four agents start training.

Performance of the agents during training is evaluated every 100,000 periods by computing the average cost per period on a fixed sample path of 100,000 periods. When costs do not improve during 30 consecutive evaluations, the agent stops training. This is equivalent to a simulation of 3 million periods, with an intermediate cost evaluation every 100,000 periods. Once all agents are finished, a new set of hyperparameters is chosen. The evaluation of one set of hyperparameters took us about 24 CPU hours. Training all experiments reported in this paper required about 3000 of these “runs”. This resulted in total training costs of around \$2500 at current Google Cloud Platform rates of \$0.033174 per virtual CPU hour. To find a “good” set of hyperparameters, investment in computing hardware is clearly indispensable. Alternatively, we relied on cloud computing.

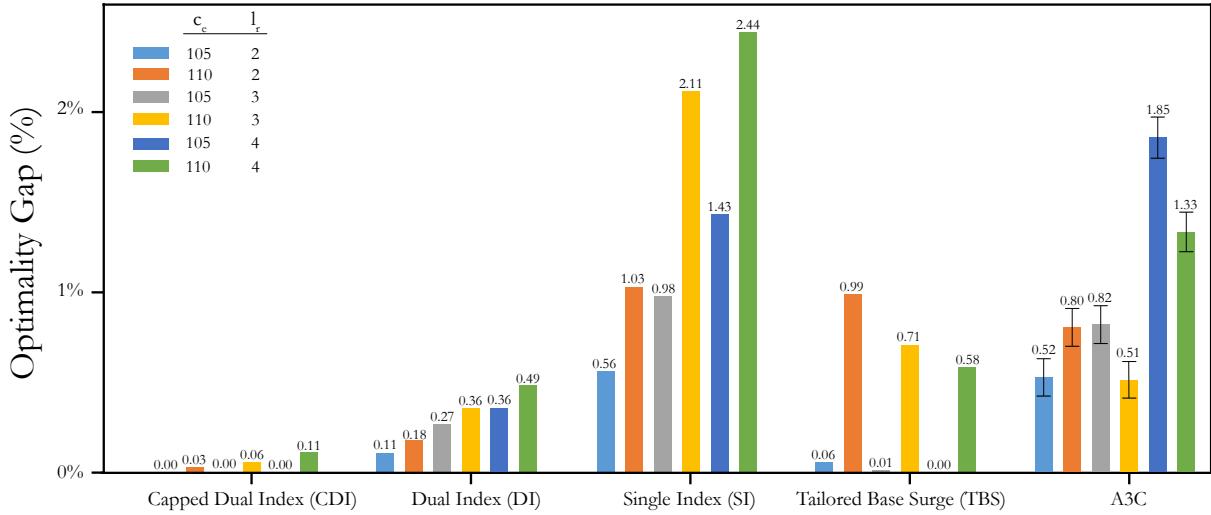


Figure 2 The A3C algorithm is able to develop policies whose expected cost is within 2% of the optimal cost. This performance is comparable to Tailored Base-Surge, Single Index and Dual Index policies, but uniformly dominated by the Capped Dual Index policy (as expected given its robust optimality). Single and Dual Index performance deteriorates as the lead time increases while TBS improves (as expected given its asymptotic optimality).

5. Performance Evaluation of DRL in Dual Sourcing Inventory Models

We evaluate the performance of DRL in the six small-scale dual sourcing settings of Veeraraghavan and Scheller-Wolf [2008] with linear sourcing costs for which the dynamic program is numerically solvable. The results are compared versus well-known heuristics. Since the value function of the dual sourcing is $L^\frac{1}{2}$ -convex, we also benchmark against the LP-ADP technique, used by Chen and Yang [2019] in a dual sourcing model with supply uncertainty.

To keep the dynamic program computationally solvable, the lead time of the expedited source is zero with a discrete uniform demand over $\{0, 1, 2, 3, 4\}$. Unit holding and backlog costs in all six experiments are $h = 5$ and $b = 495$ while the unit sourcing cost of the regular source $c_r = 100$. The six experiments range in their unit expediting sourcing cost c_e (high = 110 or low = 105) and the regular source's lead time l_r (small = 2, medium = 3 and high = 4). The optimal cost for each scenario was determined by computationally solving the DP using a discount factor of 99.9%. The costs obtained by the A3C algorithm are found by simulation and are shown with 95% confidence intervals. The other policies' parameters are optimized by simulation; after which the steady state distribution of their Markov chain was computed to compute its expected cost.

Figure 2 reports our results and shows that the A3C algorithm succeeds in developing a policy that performs within 2% of optimality in all six experiments. This performance is comparable to the Single Index and Tailored Base Surge policies. The Capped Dual Index policy, which is robustly optimal, also outperforms all other policies in this stochastic setting.

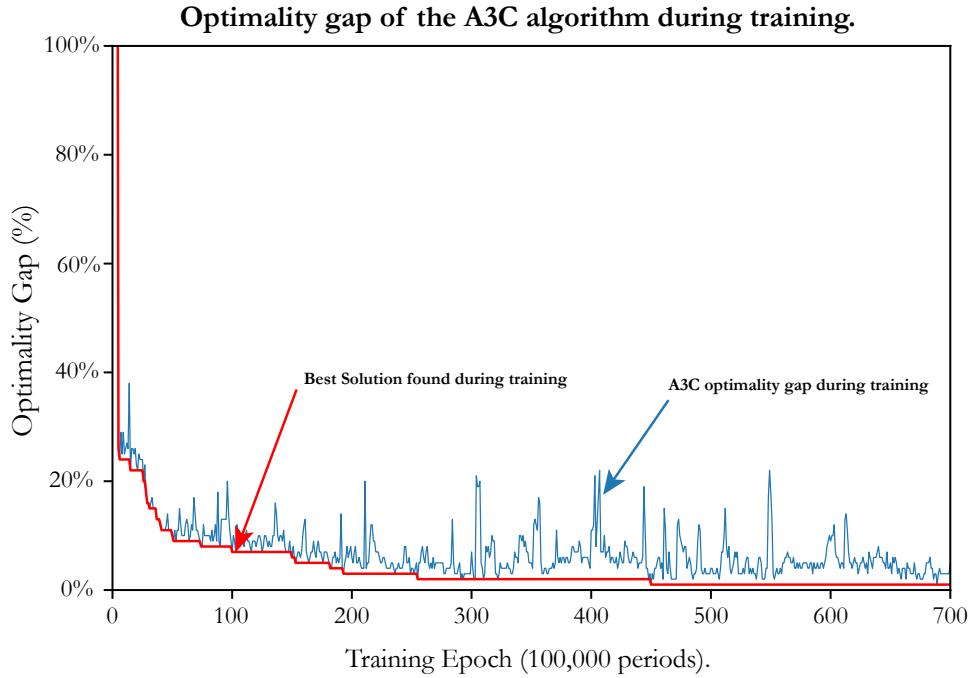


Figure 3 The A3C algorithm quickly converges to a decent policy (within $\approx 20\%$ of optimality). Then it trains fast (≈ 200 episodes of length 100,000 periods) to around 5%, after which it steadily reduces the optimality gap to less than 2% while exploring new policies that fall within 10-20% of the optimal policy.

We also compare against the LP-ADP algorithm that Chen and Yang [2019] used in a dual sourcing setting with supply uncertainty. We used the SI, DI, CDI and TBS policies as sample test policies, and for each sample path we solved the LP and evaluated the LP-greedy policies. We repeated this procedure for 100 different sample paths and tested sample paths of lengths 6,000, 11,000 and 60,000. Unfortunately, LP-ADP using this procedure did not manage to develop policies for dual sourcing that outperformed the test policies. Rather, the LP-ADP policies were often too myopic in nature; e.g., many policies did not order, or placed only expedited orders.

A potential reason may be that, in contrast to Chen and Yang [2019], our model contains no supply uncertainty and perhaps the latter is needed for LP-ADP to provide improvement over the test policies? Or perhaps the LP-ADP improvement documented by Chen and Yang [2019] in a setting with supply uncertainty resulted from using test policies that were not developed for a dual-sourcing setting with supply uncertainty? Answering these questions satisfactorily is beyond the scope of this paper. We conclude that, in our experience, the LP-ADP method was highly dependent on the sample test policy and failed to improve on the benchmark dual-sourcing heuristics that we used as sample test policies.

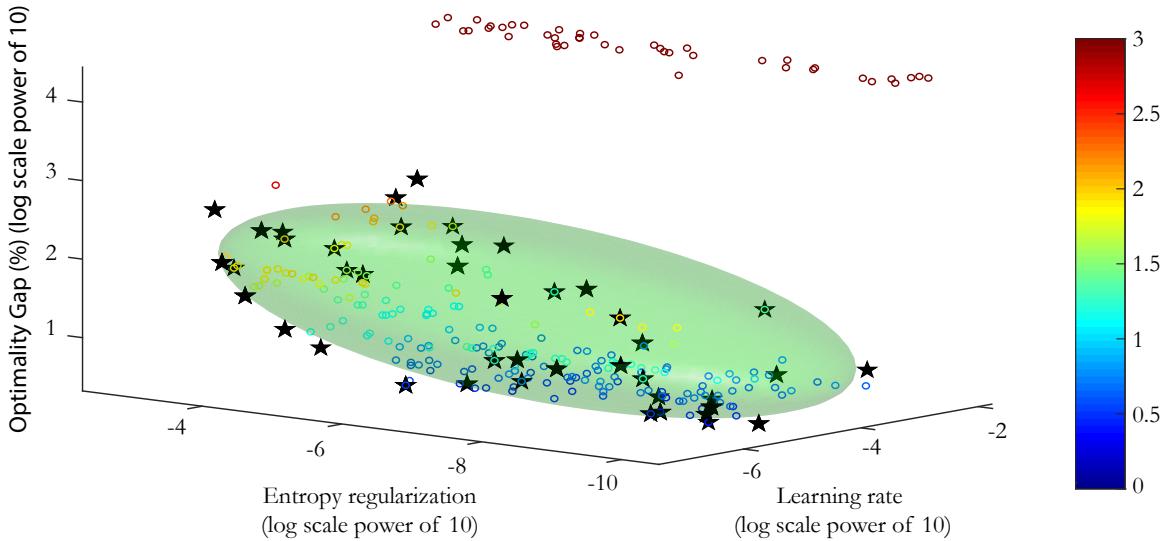


Figure 4 The automatic tuning of the learning rate and the entropy over the intervals specified in Table 1 reveals that all converging points on the convex hull (marked with black star) fit an ellipsoid (results shown for $c_e = 105$ and $l_r = 4$). The red dots indicate policies where the A3C algorithm does not develop a good policy, e.g., policies that never place orders, or policies that result in infinite inventories.

Noteworthy is that the training time of A3C remains relatively stable independent on the slow source's lead time, in contrast to dynamic programming. Nonetheless, it is not exceptional that more than five million periods are simulated per agent to find near-optimal policies, even in small-scale experiments. Figure 3 plots an example of a training curve of the A3C algorithm, where the best policy is obtained after 4.5 million periods of simulation per agent. Typically, the algorithm converges quickly to a relatively good policy (within $\approx 20\%$ of optimal cost). Then it trains fast (≈ 200 episodes of length 100,000 periods) to approximately 5%, after which it steadily reduces the optimality gap to less than 2%. These results were obtained by occasionally restarting trained models with different hyperparameter values. For instance, a smaller learning rate could occasionally improve performance of these pre-trained models.

The automatic tuning process revealed that the A3C performance was fairly insensitive to the length of the buffer size. The set of best performing values of the learning rate and entropy can be enveloped by an ellipsoid (see Figure 4). The best learning rate seems to be in the neighbourhood of 10^{-4} while the entropy factor should be set smaller than 10^{-7} to achieve good performance.

6. Performance Evaluation of DRL in Lost Sales Inventory Models

Problem setting and Formulation In a single sourcing lost sales model with lead time l , the state of the Markov Decision Process consists of the inventory on hand I_{t-1} and all outstanding

Hyperparameter	Tuning Range	Setting 1 ($p=4, l=2$)	Setting 2 ($p=4, l=3$)	Setting 3 ($p=4, l=4$)	Setting 4 ($p=9, l=2$)	Setting 5 ($p=9, l=3$)	Setting 6 ($p=9, l=4$)
Learning rate (α)	$[1 \times 10^{-5}, 1 \times 10^{-3}]$	4.26×10^{-4}	3.04×10^{-4}	3.00×10^{-4}	4.23×10^{-5}	8.8×10^{-5}	5.25×10^{-5}
Entropy regularization (β_E)	$[1 \times 10^{-5}, 1 \times 10^{-10}]$	4.74×10^{-8}	1.43×10^{-7}	2.10×10^{-7}	2.77×10^{-9}	5.43×10^{-9}	4.19×10^{-10}
Length of the Episode buffer (m)	$[1, \dots, 100]$	197	107	46	62	87	63
Action Space		$[0, 1, \dots, 20]$					

Table 2 Hyperparameters and action space of the A3C algorithm used in our six lost sales settings.

orders $Q_{t-1} = (q_{t-l}, \dots, q_{t-1})$. The state in period t , $s_t = (I_{t-1}, Q_{t-1})$ thus is l -dimensional⁵. The action a_t in period t is the order quantity q_t , so that the action space is uni-dimensional. The usual costs apply: a holding cost h per unit positive inventory, shortage costs p per unit of lost demand and ordering costs c per unit ordered. The per-period cost thus is:

$$c_t(s_t, a_t) = q_t c + h[I_t]^+ + b[I_{t-1} + q_{t-l} - d_t]^+, \quad (2)$$

in which d_t represents the stochastic demand in each period. The inventory evolves as $I_t = [I_{t-1} + q_{t-l} - d_t]^+$ while the outstanding orders shift one period: $Q_t = (q_{t-l+1}, \dots, q_t)$.

Results of Applying DRL To apply the A3C algorithm to the lost sales inventory model we modified the first input layer to reflect the l -dimensional state vector and one-dimensional action space.

We adopt six of the numerical experiments of Zipkin [2008a]: all experiments have a unit holding cost $h = 1$, ordering cost $c = 0$ and demand is Poisson distributed with $\lambda = 5$. Lead times are 2, 3 and 4 periods and the shortage penalty is either 4 or 9, resulting in six settings. The hyperparameters are automatically tuned as described in Section 4. The tuning ranges, the best performing hyperparameters and the action space we used in each experiment are reported in Table 2.

We benchmark the A3C algorithm against six heuristic policies. Three policies are simple to implement and possess interesting asymptotic behaviour: a base-stock policy, a constant-order policy and a capped base-stock policy. In addition we add three ADP policies: the myopic 1-period and myopic 2-period policies that have been shown by Zipkin [2008a] to perform well across a variety of settings, as well as the LP-ADP approach of Sun et al. [2016].

Figure 5 plots the results of our numerical experiment. The A3C algorithm performs in line with the myopic 1-period policy but cannot match the performance of the myopic 2-period policy. It does perform better than the base-stock and constant order policy but its generalization, i.e. the capped base-stock policy still performs better. We can thus conclude that the A3C algorithm, with limited modification, also develops good policies for the lost sales inventory problem.

⁵ As the outstanding orders q_{t-l} are added to the inventory prior to computing the costs, they are not taken into when computing the state space dimension.

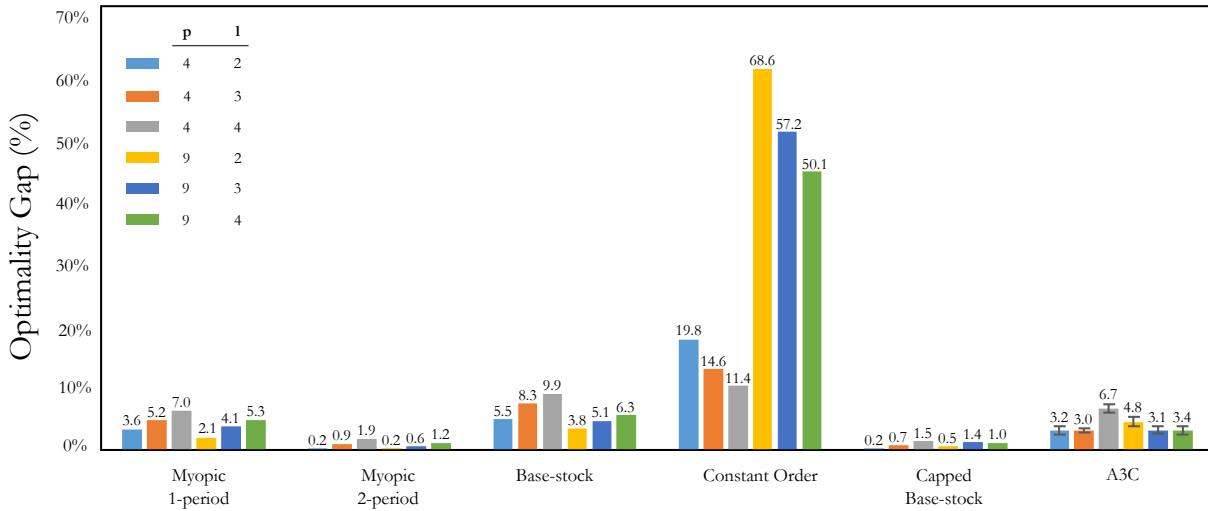


Figure 5 The A3C algorithm performs in line with the Myopic 1-period policy but cannot match the performance of the Myopic 2-period policy. It does perform better than the Base-stock and Constant Order policy but its generalization, i.e. the Capped Base-Stock policy still performs better.

Figure 6 reports the performance of 100 runs of the LP-ADP approach of Sun et al. [2016] with $p = 4$ and $l = 4$. One run corresponds to sampling 6,000 demands. Then we use the myopic policy and store the 6,000 states visited and actions taken to construct the LP. We exclude the results of the first 1,000 periods warm-up. Then we solve the LP to obtain the coefficients of the quadratic approximation. Finally, we evaluate the performance of the LP-greedy policy for 60,000 periods, excluding a warm-up of 10,000 periods. This concludes one LP-ADP run.

It is striking in Figure 6 that the results are very run dependent: only a minority of the runs achieves reasonable performance yet no run outperformed the myopic sampling test policy. Sun et al. [2016] do report occasional improvements over the sampling test policy. While LP-ADP has been reported to obtain results in line or slightly better than the sampling test policy, in our experience it does not appear to be a stable, nor well-performing method for our lost sales inventory problem.⁶ The benefit of LP-ADP is that it is computationally less intensive than A3C. One LP-ADP run can typically be finished in minutes instead of hours. Even 100 runs was easily done on one personal computer without requiring cloud computing.

The automatic tuning of the hyper-parameters for the lost sales inventory model reveals similarities with dual sourcing. The A3C performance is fairly insensitive to the length of the buffer size and the set of best performing values of the learning rate and entropy can be enveloped by an ellipsoid (see Figure 7). The best learning rate seems to be in the neighbourhood of 10^{-4} while the entropy factor should be set smaller than 10^{-7} to achieve good performance.

⁶ We tested different lengths of sample paths but did not observe any improvement in performance.

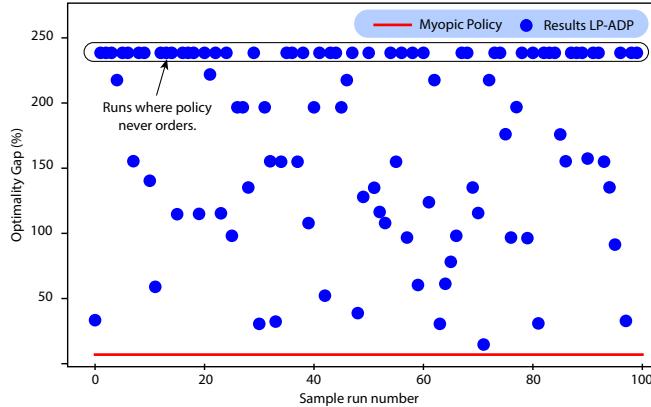


Figure 6 Sample run performance of the LP-ADP method relative to the myopic base policy for a lost sales problem with $p = 4$ and $l = 4$.

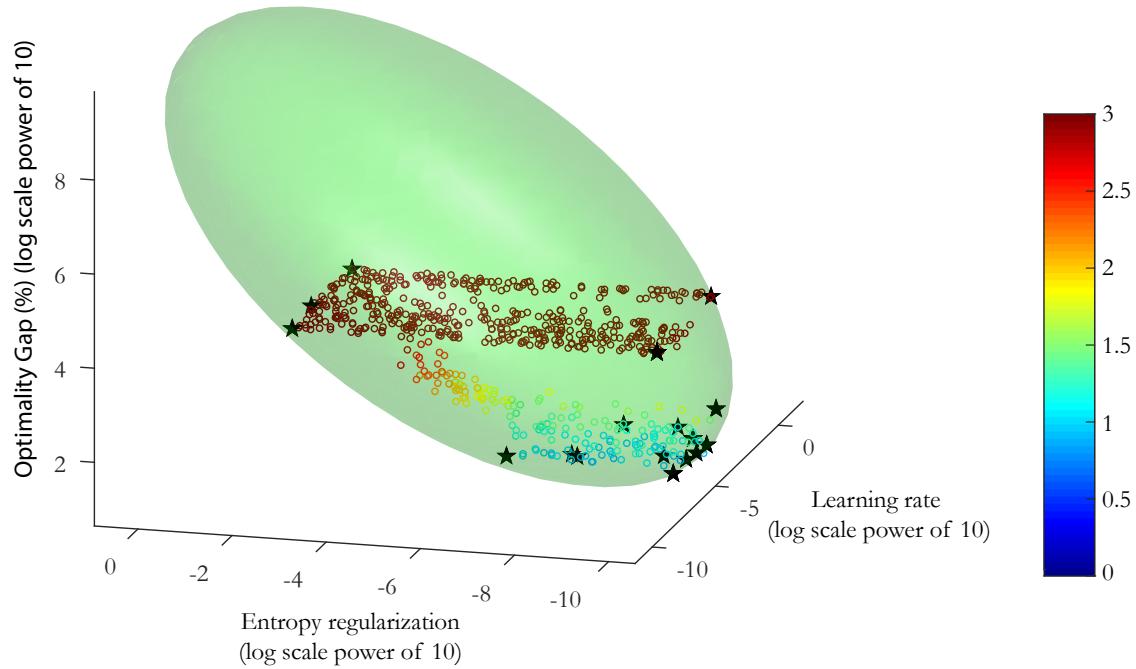


Figure 7 The automatic tuning of the learning rate and the entropy over the intervals specified in Table 1 reveals that all points on the convex hull (marked with black star) fit an ellipsoid (results shown for $p = 4$ and $l = 4$). The red dots indicate policies where the A3C algorithm does not develop a good policy, e.g., policies that never place orders, or policies that result in infinite inventories.

7. Performance Evaluation of DRL in Multi-echelon Inventory Models

Problem setting and Formulation We adopt the multi-echelon inventory model of Van Roy et al. [1997] with one warehouse and K identical retailers (see Figure 8). Van Roy et al. [1997] are the first to adopt a neuro-dynamic programming approach in inventory management; They use a neural net to approximate the value function which is updated using temporal difference learning.

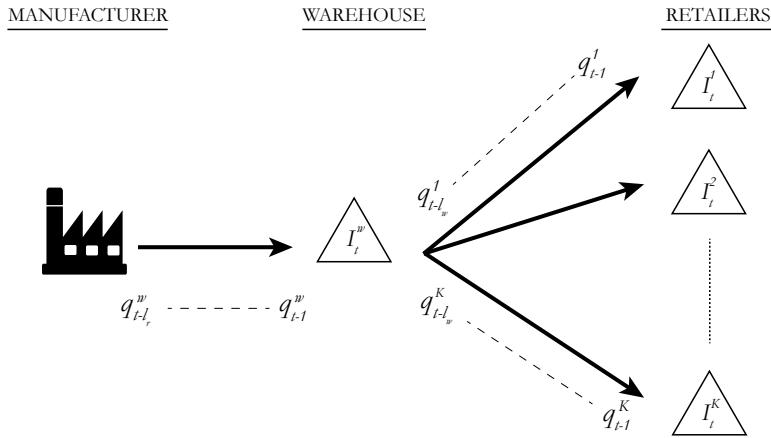


Figure 8 Two-echelon supply chain of Van Roy et al. [1997] in which one capacitated warehouse is replenished from a manufacturer with lead time l_r and replenishes K capacitated retailers with lead time l_w .

Inventory held at the central warehouse or at the retailers incurs a unit holding cost h_w and h_r , respectively. Transshipment of inventory between retailers is not allowed. The replenishment lead time from the manufacturer to the warehouse is l_w periods, and from the warehouse to any retailer l_r periods. All locations have limited capacity: (1) maximum production rate is C^m units per period; (2) warehouse inventory position (incl. outstanding orders) cannot exceed C^w units; and (3) retail inventory position (incl. outstanding orders) cannot exceed C^r units.

The model's MDP has state $s_t = (I_{t-1}^w, Q_{t-1}^w, I_{t-1}^r, Q_{t-1}^r)$, where I_{t-1}^w is the end-of-period inventory at the warehouse and $Q_{t-1}^w = (q_{t-l_w}^w, \dots, q_{t-1}^w)$ its outstanding orders; $I_{t-1}^r = (I_{t-1}^1, \dots, I_{t-1}^K)$ the end-of-period inventories at the K retailers, and $Q_{t-1}^r = (q_{t-l_r}^1, \dots, q_{t-1}^1, q_{t-l_r}^2, \dots, q_{t-1}^2, \dots, q_{t-l_r}^K, \dots, q_{t-1}^K)$ their respective outstanding orders. The state space is thus $(l_w + Kl_r)$ -dimensional⁷.

The action vector $a_t = (q_t^w, q_t^1, \dots, q_t^K)$, consists of the order quantities: q_t^w at the warehouse and q_t^i ($i = 1, \dots, K$) at retailer i . The retailer order quantities cannot exceed the end-of-period inventory on hand in the warehouse, i.e. $\sum_{i=1}^K q_t^i < I_{t-1}^w + q_{t-l_w}^w$. The sequence of events is as follows. Each period t demand d_i at retailer i is fulfilled by the inventory on hand at retailer i . In case of a stock-out, $[d_i - (I_{t-1}^i + q_{t-l_r}^i)]^+ > 0$, customers either decide to wait for an expedited, same-day delivery from the warehouse — this occurs with probability P_w — or walk away such that the retailer incurs a lost sale with probability $1 - P_w$. The system is thus a hybrid of a backlogging and lost sales model. Important is that expedited deliveries are shipped *after* retailers have placed orders (q_t^1, \dots, q_t^K) and are allocated inventory from the warehouse. This implies that only

⁷ Note that Q_{t-1}^w and Q_{t-1}^r have a length l_w and $K \times l_r$ but only the first $(l_w - 1)$ and $(l_r - 1)$ elements of the pipeline vectors are added to the state-dimension as $q_{t-l_w}^w$ and $q_{t-l_r}^i$ are added to the inventories prior to incurring the cost. Together with $K + 1$ inventories the state dimension is thus $(l_w + Kl_r)$.

$[I_{t-1}^w + q_{t-l_w}^w - \sum_{i=1}^K q_t^i]$ units can be expedited. Let B_t denote the number of customers in period t that request an expedited delivery. The number of requested expedite deliveries exceeding the available end-of-period inventory on hand at the warehouse, $\left[B_t - [I_{t-1}^w + q_{t-l_w}^w - \sum_{i=1}^K q_t^i] \right]^+$, is lost and the remainder, $B_t - \left[B_t - [I_{t-1}^w + q_{t-l_w}^w - \sum_{i=1}^K q_t^i] \right]^+$, is expedited and shipped same-day from the warehouse. Expedited deliveries by the warehouse incur a unit cost c_w while unmet demand is lost at a unit penalty cost p . The cost in period t thus is:

$$\begin{aligned} c_t(s_t, a_t) = & h_w I_t^w + h_r \sum_{i=1}^K I_t^i \\ & + c_w \left[B_t - \left[B_t - [I_{t-1}^w + q_{t-l_w}^w - \sum_{i=1}^K q_t^i] \right]^+ \right] \\ & + p \left[\left[B_t - [I_{t-1}^w + q_{t-l_w}^w - \sum_{i=1}^K q_t^i] \right]^+ + \sum_{i=1}^K [d_i - I_{t-1}^i + q_{t-l_r}^i]^+ - B_t \right], \quad (3) \end{aligned}$$

where the first two elements represent the inventory holding costs in the supply chain, the third element the expedited delivery costs and the last two elements the lost sales costs upon stock-out at the warehouse and upon stock-out at the retailer. The state evolves as follows: $I_t^w = I_{t-1}^w + q_{t-l_w}^w - \sum_{i=1}^K q_t^i - \left[B_t - \left[B_t - [I_{t-1}^w + q_{t-l_w}^w - \sum_{i=1}^K q_t^i] \right]^+ \right]$, $I_t^i = [I_{t-1}^i + q_{t-l_r}^i - d_i]^+$ for each retailer i , $Q_t^w = (q_{t-l_w+1}^w, \dots, q_t^w)$ and $Q_t^i = (q_{t-l_r+1}^i, \dots, q_t^i)$ for each retailer i .

Results of Applying DRL To apply A3C in the multi-echelon inventory model, we adapt the width of the first layer of the neural net to the dimension of the state space ($l_w + K l_r$). As the action space is K -dimensional, we reduce the action space by adopting a base-stock replenishment policy for the warehouse and the retailers with state-dependent order-up-to levels, similar to Van Roy et al. [1997]. This reduces the action space to only two dimensions $a_t = [S_{s_t}^w, S_{s_t}^r]$, consisting of a base-stock level at the warehouse and a base-stock level for all retailers, which all depend on the state s_t . The order quantities are then determined as follows: The warehouse order q_t^w raises its inventory position: $q_t^w = S_{s_t}^w - (I_{t-1}^w + \sum_{i=1}^{l_w} q_{t-i}^w)$ to its base-stock level $S_{s_t}^w$. Note that q_t^w may not exceed production capacity C^w and the inventory position may not exceed C^w . Hence,

$$q_t^w = \min \left(C^w - (I_{t-1}^w + \sum_{i=1}^{l_w} q_{t-i}^w), \min \left(S_{s_t}^w - (I_{t-1}^w + \sum_{i=1}^{l_w} q_{t-i}^w), C^w \right) \right).$$

After the warehouse has ordered, each retailer places its order to raise their inventory position to the base-stock level $S_{s_t}^r$: $q_t^i = S_{s_t}^r - (I_{t-1}^i + \sum_{i=1}^{l_r} q_{t-i}^i)$. After ordering, each retailer's inventory position may not exceed C_r :

$$q_t^i \leq C_r - (I_{t-1}^i + \sum_{i=1}^{l_r} q_{t-i}^i).$$

Hyperparameters	Tuning Range	Setting 1	Setting 2
Learning rate used	$[10^{-5}, 10^{-3}]$	5.78×10^{-5}	1.74×10^{-4}
Entropy regularization (β_E)	$[10^{-10}, 10^{-6}]$	4.68×10^{-5}	1.46×10^{-8}
Length of the Episode buffer (m)	100	100	100
Base-stock levels warehouse (S_w)		$[50, 60, \dots, 100]$	$[50, 60, \dots, 100]$
Base-stock levels retailers (S_r)		$[0, 5, \dots, 40]$	$[0, 5, \dots, 50]$

Table 3 Hyperparameters of the A3C algorithm to achieve results in both multi-echelon settings.

If the warehouse has insufficient capacity to deliver all retailer orders, the inventory is allocated to the retailers by minimizing the difference in inventory position between all retailers after orders are allocated following Van Roy et al. [1997].

We test the two numerical settings in Van Roy et al. [1997] with 10 retailers and varying values of lead times and demand parameters (see Figure 9). Each element of the random demand vector $D = (d_1, \dots, d_K)$ is sampled from a normal distribution with mean μ and standard deviation σ , rounded to the nearest integer and truncated at zero: $d_i = [\mathcal{N}(\mu, \sigma) - \frac{1}{2}]^+$.

As computing the optimal policy is intractable we only compare the A3C versus a base-stock policy with constant, i.e., not state-dependent, base-stock levels S_w and S_r . Unlike the dual sourcing and lost sales inventory models we cannot deploy the LP-ADP approach, as we do not know the structure of the value function, and the K -dimensional demand makes the number of transition probabilities too large and the formulation of the constraints of the LP infeasible.

We trained the A3C algorithm using the automatic tuning process with the same set of hyperparameters as in dual sourcing and lost sales (see Table 3). We fixed the buffer size at 100, which provided decent results after 100 training runs of the A3C algorithm. The discrete action space of the A3C algorithm is identical to that used in Van Roy et al. [1997]. The state-dependent base-stock levels $S_{s_t}^w$ are restricted to 6 values $\{50, 60, \dots, 100\}$ while $S_{s_t}^r$ is restricted to 9 values $\{0, 5, 10, \dots, 40\}$ in Setting 1 and 11 values $\{0, 5, 10, \dots, 50\}$ in Setting 2. Figure 9 shows that the A3C algorithm improves 8.95% and 12.09% on the base-stock policy with constant base-stock levels. These results are in line with, and slightly better in setting 2, than Van Roy et al. [1997] who report savings around 10%.

One key distinguishing factor with the method of Van Roy et al. [1997] is that A3C (or DRL) does not require manual feature engineering. Van Roy et al. [1997] manually develop 23 features based on the state vector. Instead, DRL or A3C learns directly from the state vector. This is possible because our network uses four layers and is ‘deeper’ than the single layer with activation neural net used in Van Roy et al. [1997]. Their careful selection of the features has resulted in excellent performance but required problem knowledge. In contrast, DRL does not require any manual feature engineering, a desirable characteristic of a general purpose technology. Moreover, we observe that A3C manages to find slightly better policies in the larger setting 2.

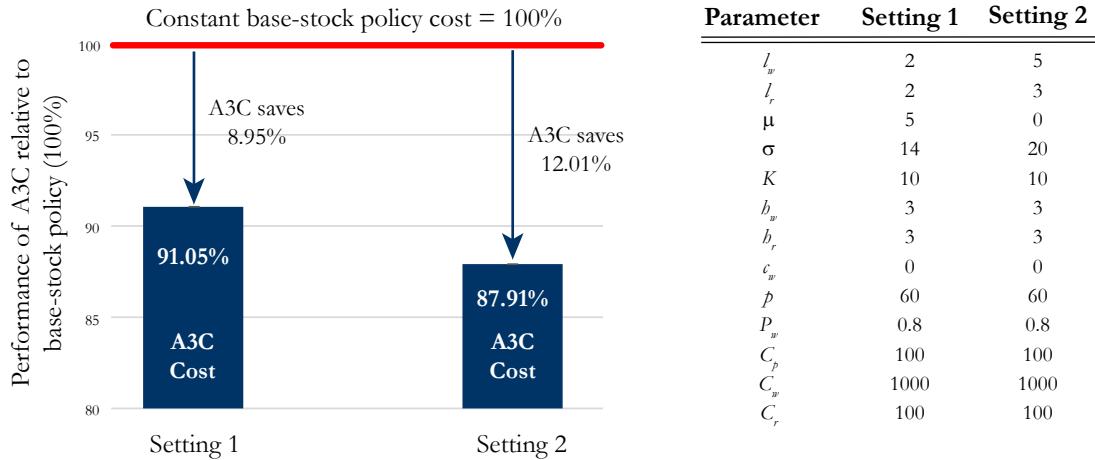


Figure 9 In both settings A3C significantly outperforms the base-stock policy with $8.95 \pm 0.13\%$ and $12.09 \pm 0.39\%$ (where the errors indicate 95% confidence) for Settings 1 and 2, respectively. Error bars at $\pm 1.96\sigma$ indicate 95% confidence interval. These results are in line (and slightly better in the larger setting) than Van Roy et al. [1997] who report savings around 10%. Moreover, A3C does not use any feature engineering and learns directly from the state space.

8. Conclusions and Reflections

Our study provides evidence that deep reinforcement learning can effectively solve classic, intractable inventory problems. The A3C algorithm can be trained to produce policies that match the performance of state-of-the-art heuristics and other approximate dynamic programming methods. After extensive manual tuning of 9 hyper-parameters in one problem setting, we found that fixing 6 parameters and automatically tuning the learning rate, entropy regularization, and buffer size resulted in good performance in the other two problem settings. With such minimal changes among three different problem settings, DRL thus seems a promising general purpose technology.

Yet we also found that A3C does not outperform all policies and remains, like other ADP methods, a black box. Tuning the hyper-parameters of the neural network is effort-intensive and requires both art and science. The computational requirements are formidable: despite the rapid improvement in computing technology, we had to resort to cloud computing.

So our conclusion is nuanced: DRL seems a promising general purpose technology that can be applied to intractable inventory problems. Yet, there remains a job for academic inventory researchers to generate structural policy insight or design specialized policies that are (ideally provably) near optimal.

Our study has focused on solving classic, stationary models of inventory management. Before applying DRL in practice, we should evaluate performance in non-stationary environments with multiple products and unknown demands that must be learned from historical data. We hope that

by sharing our experience and code, this paper will provide a stepping stone for such research in inventory management and other operational settings.

References

- Allon G, Van Mieghem JA (2010) Global Dual Sourcing: Tailored Base-Surge Allocation to Near- and Offshore Production. *Mgt Sci* 56(1):110–124.
- Arrow KJ, Karlin S (1958) *Studies in the mathematical theory of inventory and production* (Stanford University).
- Bellman R (1954) The Theory of Dynamic Programming. *Bulletin of the Amer Math Soc* 60(6):503–515.
- Brown DB, Smith JE (2014) Information relaxations, duality, and convex stochastic dynamic programs. *Ops Res* 62(6):1394–1415.
- Chen W, Dawande M, Janakiraman G (2014) Fixed-dimensional stochastic dynamic programs: An approximation scheme and an inventory application. *Ops Res* 62(1):81–103.
- Chen W, Yang H (2019) A heuristic based on quadratic approximation for dual sourcing problem with general lead times and supply capacity uncertainty. *IIE Transactions*.
- Clark AJ, Scarf H (1960) Optimal Policies for a Multi-echelon Inventory Problem. *Mgt Sci* 6(4):363–505.
- de Kok T, Grob C, Laumanns M, Minner S, Rambau J, Schade K (2018) A typology and literature review on stochastic multi-echelon inventory models. *Eur J of Op Res* 269(3):955–983.
- Fang J, Zhao L, Fransoo JC, Van Woensel T (2013) Sourcing strategies in supply risk management: An approximate dynamic programming approach. *Comput. Oper. Res.* 40(5):1371–1382.
- Federgruen A, Zipkin P (1984) Approximations of dynamic, multilocation production and inventory problems. *Mgt Sci* 30(1):69–84.
- Fukuda Y (1964) Optimal Policies for the Inventory Problem with Negotiable Leadtime. *Mgt Sci* 10(4):690–708.
- Giannoccaro I, Pontrandolfo P (2002) Inventory management in supply chains: a reinforcement learning approach. *Int. Journ. Prod. Econ.* 78(2):153–161.
- Goldberg DA, Katz-Rogozhnikov DA, Lu Y, Sharma M, Squillante MS (2016) Asymptotic Optimality of Constant-Order Policies for Lost Sales Inventory Models with Large Lead Times. *Mathematics of Operations Research* 41(3):898–913.
- Halman N, Klabjan D, Mostagir M, Orlin J, Simchi-Levi D (2009) A fully polynomial-time approximation scheme for single-item stochastic inventory control with discrete demand. *Math of Ops Res* 34(3):674–685.
- Hua Z, Yu Y, Zhang W, Xu X (2015) Structural properties of the optimal policy for dual-sourcing systems with general lead times. *IIE Transactions* 47(8):841–850.

- Huh WT, Janakiraman G, Muckstadt JA, Rusmevichientong P (2009) Optimality of Order-Up-To Policies in Lost Sales Inventory Systems. *Mgt Sci* 55(3):404–420.
- Keller PW, Mannor S, Precup D (2006) Automatic basis function construction for approximate dynamic programming and reinforcement learning. *Proceedings of the 23rd International Conference on Machine Learning*, 449–456, ICML '06 (ACM).
- Kingma DP, Ba J (2015) Adam: A Method for Stochastic Optimization. *ICLR 2015*, URL <http://arxiv.org/abs/1412.6980>.
- Kohl N, Stone P (2004) Policy gradient reinforcement learning for fast quadrupedal locomotion. *Robotics and Automation, 2004. Proceedings. ICRA'04. 2004 IEEE International Conference on*, volume 3, 2619–2624 (IEEE).
- Mnih V, Badia AP, Mirza M, Graves A, Lillicrap TP, Harley T, Silver D, Kavukcuoglu K (2016) Asynchronous Methods for Deep Reinforcement Learning. *arXiv preprint arXiv:1602.01783* 48.
- Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, Bellemare MG, Graves A, Riedmiller M, Fidjeland AK, Ostrovski G, Petersen S, Beattie C, Sadik A, Antonoglou I, King H, Kumaran D, DaanWierstra, Legg S, Hassabis D (2015) Human-level control through deep reinforcement learning. *Nature* 518:529–533.
- Morton TE (1969) Bounds on the Solution of the Lagged Optimal Inventory Equation with no Demand-Backlogging and Proportional Costs. *SIAM Rev.* 11(4):572–596.
- Morton TE (1971) The Near-Myopic Nature of the Lagged-Proportional-Cost Inventory Problem with Lost-Sales. *Ops Res* 19(7):1708–1716.
- Murota K (2003) *Discrete Convex Analysis: Monographs on Discrete Mathematics and Applications* 10 (Society for Industrial and Applied Mathematics).
- Nahmias S, Smith SA (1993) *Sarin R.K. (eds) Perspectives in Operations Management* (Springer, Boston, MA).
- Nahmias S, Smith SA (1994) Optimizing Inventory Levels in a Two Echelon Retailer System with Partial Lost Sales. *Ops Res* 40(5):582–596.
- Oroojlooyjadid A, Nazari M, Snyder L, Takáč M (2017) A Deep Q-Network for the Beer Game with Partial Information. *arXiv preprint arXiv:1708.05924* .
- Oroojlooyjadid A, Snyder L, Takáč M (2016) Applying Deep Learning to the Newsvendor Problem. *arXiv preprint arXiv:1607.02177* .
- Powell WB (2011) *Approximate dynamic programming: solving the curses of dimensionality*, volume 6788 (Wiley).
- Puterman ML (1994) *Markov Decision Processes: Discrete Stochastic Dynamic Programming* (New York, NY, USA: John Wiley & Sons, Inc.), 1st edition.

- Reiman MI (2004) A new and simple policy for the continuous review lost sales inventory model. *Working Paper, Bell Labs, Lucent Technologies* .
- Scheller-Wolf A, Veeraraghavan S, van Houtum GJ (2003) Effective dual sourcing with a single index policy. *Working Paper, Wharton at the University of Pennsylvania* .
- Schulman J, Moritz P, Levine S, Jordan M, Abbeel P (2015) High-Dimensional Continuous Control Using Generalized Advantage Estimation. *arXiv preprint arXiv:1506.02438* .
- Sheopuri A, Janakiraman G, Seshadri S (2010) New Policies for the Stochastic Inventory Control Problem with Two Supply Sources. *Ops Res* 58(3):734–745.
- Silver D, Schrittwieser J, Simonyan K, Antonoglou I, Huang A, Guez A, Hubert T, Baker L, Lai M, Bolton A, Chen Y, Lillicrap T, Hui F, Sifre L, van den Driessche G, Graepel T, Hassabis D (2017) Mastering the game of Go without human knowledge. *Nature* 550(7676):354–359.
- Strehl AL, Li L, Wiewiora E, Langford J, Littman ML (2006) Pac model-free reinforcement learning. *Proceedings of the 23rd international conference on Machine learning*, 881–888 (ACM).
- Sun J, Van Mieghem JA (2018) Robust Dual Sourcing Inventory Management: Optimality of Capped Dual Index Policies. *M&SOM* Forthcoming.
- Sun P, Wang K, Zipkin P (2016) Quadratic Approximation of Cost Functions in Lost Sales and Perishable Inventory Control Problems. *Working Paper, Fuqua School of Business, Duke University, Durham, NC 27708, USA* .
- Sutton RS, Barto AG (1998) *Introduction to reinforcement learning*, volume 135 (MIT press Cambridge).
- Tesauro G (1995) Temporal difference learning and td-gammon. *Communications of the ACM* 38(3):58–68.
- Van Roy B, Bertsekas DP, Lee Y, Tsitsiklis JN (1997) A Neuro-Dynamic Programming Approach to Retailer Inventory Management. *Proceedings of the IEEE Conference on Decision and Control* .
- Veeraraghavan S, Scheller-Wolf A (2008) Now or Later: A Simple Policy for Effective Dual Sourcing in Capacitated Systems. *Ops Res* 56(4):850–864.
- Watkins C (1989) *Learning from delayed rewards*. Ph.D. thesis, Cambridge University.
- Whittemore AS, Saunders SC (1977) Optimal Inventory Under Stochastic Demand With Two Supply Options. *J on Appl Math* 32(2):293–305.
- Xin L (2019) Understanding the performance of capped base-stock policies in lost-sales inventory models. *Working Paper, The University of Chicago Booth School of Business* .
- Xin L, Goldberg DA (2017) Asymptotic Optimality of Tailored Base-Surge Policies in Dual-Sourcing Inventory Systems. *Mgt Sci* 64(1):437–452.
- Zipkin P (2008a) Old and new methods for lost-sales inventory systems. *Ops Res* 56(5):1256–1263.
- Zipkin P (2008b) On the Structure of Lost-Sales Inventory Models. *Ops Res* 56(4):937–934.