

# The Title of your Report

Anonymous Authors

**Abstract**—A short summary of your project. You should change also the title, but do *not* enter any author names or anything that unnecessarily identifies any of the authors. It is suggested you use a similar structure (sections, etc.) as demonstrated in this document, but you can make the section headings more descriptive if you wish. Of course you should delete all the text in this template and write your own! – this text simply provides detailed instructions/hints on how to proceed.

## I. INTRODUCTION

In the following paper we implement three different reinforcement learning agents and evaluate their performance within a supply chain environment. All the code used to create the results presented in this paper can be found here <sup>1</sup>.

The supply chain optimization environment represents a problem faced by companies whose supply chain consists of a factory and multiple warehouses (so called hub-and-spoke networks as described by [4]). The main decision that needs to be made in these settings is how many products should be produced in the factory and how much stock should be build up in the warehouses. Seasonal demand can further complicate the decision problem since it might require the companies to start building up stock early so that high demands in the future can be satisfied (e.g. think about christmas where the stock of eggnog needs to be build up during november and december to be able to satisfy the high demand levels during the holidays). This requirement for mid- to long-term planning indicates that the problem can be formulated as a multi-step decision problem and further exhibits a weak signal because the reward for an early build-up of stock will only be realized in the future when the actual demand reaches its maximum. Even for small supply chain networks we find that due to the curses of dimensionality, as described by [2], the state- and action-spaces of the respective decision problems become unfeasibly large. Therefore we turn to function-approximation and policy-search methods of reinforcement learning that are less affected by these problems. In our case we chose approximate SARSA and the REINFORCE algorithm as a basis for the agents.

## II. BACKGROUND AND RELATED WORK

The supply chain optimization problem presented in this paper will be modeled as a markov decision process (MDP). Furthermore we rely on the approximate SARSA and REINFORCE algorithms to design agents capable of finding a good strategy to navigate within the environment. MDPs are multi-step stochastic decision problems that rely on the markov property which implies that the transition probability  $P(s_{n+1} | s_n)$  between two states  $s_{n+1}$  and  $s_n$  only

relies on the current state  $s_n$ . We describe the MDP similar to [1] and [2] as an eight-tuple  $(S, D, A, \mathcal{X}, T, Q, r, \gamma)$  consisting of a state-space  $S$ , a random environment process  $D$ , an action space  $A$ , a set of feasible actions  $\mathcal{X}$ , a transition function  $T$ , transition probabilities  $Q$ , a one-step reward function  $r$  and some discount factor  $\gamma$ .

The designed agents are based on the blueprints for the approximate SARSA and REINFORCE algorithm from [3]. We compare the performance of these reinforcement learning algorithms with an agent that acts according to a fixed heuristic based on the  $(\varsigma, Q)$ -Policy <sup>2</sup> as described by [5].  $(\varsigma, Q)$ -Policies induce a simple replenishment strategy where the stock of a warehouse will be replenished by some amount  $Q$  as soon as it falls under a threshold  $\varsigma$ . A basic visualization of the  $(\varsigma, Q)$ -Policy is depicted in Figure 1.

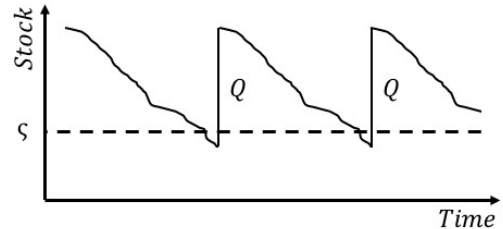


Fig. 1. Schematic visualization of a  $(\varsigma, Q)$ -Policy for a single warehouse.

Furthermore we use several different stepsize formulars from [2] that are used to update some of the agents parameters such as the learning rate  $\alpha$ . The idea of specific stepsize rules is to allow the agent to go through a longer learning phase in the beginning with higher learning rates and then converge in the long term to stepsizes close to 0. Among others we use a so called search-then-converge stepsize rule as presented in [2] where the learning rate in the  $n^{th}$  iteration is calculated by

$$\alpha_n = \alpha_0 \frac{\frac{b}{n} + a}{\frac{b}{n} + a + n^\beta} \quad (1)$$

with some parameters  $\alpha_0, b, a$  and  $\beta$  that can be set by the user based on the required behavior of the stepsize.

## III. THE ENVIRONMENT

We will model the environment as a MDP that we regard over an infinite number of periods  $t = 0, 1, \dots$ . The environment consists of one factory (indexed by 0) and up to  $k \in \mathbb{N}$  warehouses (indexed by  $j = 1, \dots, k$ ) where in each period it needs to be decided how many units of a product (e.g. butter) should be produced and how many should be shipped to the individual warehouses. A representation of a network with 5 warehouses is depicted in Figure 2. The state

<sup>1</sup>Our code is available here: <http://anonymouslinktoyourcode.zip>

<sup>2</sup>In the literature this policy is usually called  $(\varsigma, Q)$ -Policy. We chose a slightly different name to avoid notation conflicts in our model description.



Fig. 2. Example of a supply chain network with  $k = 5$  warehouses ( $s_1, \dots, s_5$ ) and the factory ( $s_0$ ).

space consists of  $k + 1$  elements describing the current stock levels of the factory and the warehouses ( $s_j, j = 0, \dots, k$ ), where each stock level is limited by some capacity  $c_j \in \mathbb{N}$ . Note that this implies that the factory itself can also build up stock. Furthermore we assume an environment process with states  $d \in D \in \{0, \dots, d_{max}\}^k, d_{max} \in \mathbb{N}$  describing the (stochastic) demand  $d_j$  at each warehouse with transition probabilities  $q_d(d')$  that are unknown by the agent. We assume a discretized version of a sinus-function and some shocks  $\epsilon$  where the demand follows a model

$$d_{i,t} = \left\lceil \frac{d_{max}}{2} \sin\left(\frac{2\pi(t+2i)}{12}\right) + \frac{d_{max}}{2} + \epsilon_{i,t} \right\rceil \quad (2)$$

with  $\lceil \cdot \rceil$  the ceiling function and  $P(\epsilon_{i,t} = 0) = P(\epsilon_{i,t} = 1) = 0.5$ . We will include  $d_{last}, d_{hist} \in D$  in our state space where for some time  $t$  we have the tuple  $\tilde{s} = (s, d_{last}, d_{hist}) \in \tilde{S} := S \times D^2$  with  $s = s_t, d_{last} = d_{t-1}$  and  $d_{hist} = d_{t-2}$ . Note that this implies that the demand  $d_t$  of period  $t$  is not observed by the agent. This means, that the demand in a period  $t$  will be realized after observing the stock levels  $s_t$  and can only be observed in the next period  $t + 1$ . The reason we add the old demands to the state space is to allow the agent to have a simple understanding of the demand history to be able to gather a basic understanding of the demands movement. In each period the agent can now set the factories production level for the next period  $a_0 \in \{0, \dots, \rho_{max}\}$  (with a maximum production of  $\rho_{max} \in \mathbb{N}$ ) as well as the number of products shipped to each location  $a_j \in \mathbb{N}^k$  that is naturally limited by the current storage level in the factory ( $\sum_{j=1}^k a_j \leq s_0$ ). Based on this information and the demand  $d$  we can now describe the state transitions  $(\tilde{s}, d, a) \rightarrow T(\tilde{s}, d, a)$  by

$$\begin{aligned} T(\tilde{s}, d, a) := & (\min\{s_0 + a_0 - \sum_{j=1}^k a_j, c_0\}, \\ & \min\{s_1 + a_1 - d_1, c_1\}, \\ & \dots \\ & \min\{s_k + a_k - d_k, c_k\}, \\ & d, \\ & d_{last}). \end{aligned} \quad (3)$$

The reward within each period consists of the revenue from sold products at a fix price  $p$  less production costs  $\kappa_{pr}a_0$ , storage costs  $\sum_{j=0}^k \kappa_{st,j} \max\{s_j, 0\}$ ,

penalty costs  $\kappa_{pe} \sum_{j=1}^k \min\{s_j, 0\}$  and transportation costs  $\sum_{j=1}^k \kappa_{tr,j} \lceil a_j / \zeta_j \rceil$ . We let  $p, \kappa_{pr}, \kappa_{pe}, \kappa_{st,j}, \kappa_{tr,j} \in \mathbb{R}_+$  and  $\zeta_j \in \mathbb{N}$ . We can now define the one-step reward function by

$$\begin{aligned} r(\tilde{s}, d, a) := & p \sum_{j=1}^k d_j - \kappa_{pr}a_0 - \sum_{j=0}^k \kappa_{st,j} \max\{s_j, 0\} \\ & + \kappa_{pe} \sum_{j=1}^k \min\{s_j, 0\} - \sum_{j=1}^k \kappa_{tr,j} \left\lceil \frac{a_j}{\zeta_j} \right\rceil. \end{aligned} \quad (4)$$

Furthermore we chose a discounting factor  $\gamma \in (0, 1)$  that can be interpreted e.g. as a result of inflation.

Based on this model description we can now formulate the task as an infinite horizon markov decision process that is subject to a varying environment. Thereby we receive the tuple  $(\tilde{S}, D, A, \mathcal{X}, T, Q, r, \gamma)$  with

#### Definition 1.

1.  $\tilde{S} \times D := \prod_{j=0}^k \{s_j \in \mathbb{Z} \mid s_j \leq c_j\} \times D^3, c_j \in \mathbb{N}$  the state space consisting of elements  $(\tilde{s}, d) = ((s, d_{last}, d_{hist}), d)$ ;
2.  $A := \{0, \dots, \rho_{max}\} \times \mathbb{N}_0^k$  the action space consisting of elements  $a = (a_0, \dots, a_k)$ ;
3.  $\mathcal{X}(\tilde{s}) := \{0, \dots, \rho_{max}\} \times \{a \in \mathbb{N}_0^k \mid \sum_{j=1}^k a_j \leq s_0\}$ , the set of all feasible actions in a state  $\tilde{s}$  and  $\mathcal{X} := \{(\tilde{s}, d, a) \in \tilde{S} \times D \times A \mid a \in \mathcal{X}(\tilde{s})\}$ ;
4.  $T : \tilde{S} \times D \times A \rightarrow \tilde{S}$  the transition function defined by Eq. (3);
5.  $Q : \mathcal{X} \times \tilde{S} \times D \rightarrow [0, 1]$  the transition probabilities with  $Q(\tilde{s}', d' \mid \tilde{s}, d, a) := q_d(d')$  for  $\tilde{s}' = T(\tilde{s}, d, a)$  and 0 otherwise;
6.  $r : \tilde{S} \times D \times A \rightarrow \mathbb{R}$  the one-step reward function as described in Eq. (4);
7.  $\gamma \in (0, 1)$  the discounting factor.

For the calculation of the expected discounted revenue for each tuple  $(\tilde{s}, d)$  we obtain the value function  $V$  with

$$\begin{aligned} V(\tilde{s}, d) = & \max_{a \in \mathcal{X}(\tilde{s})} \{r(\tilde{s}, d, a) \\ & + \gamma \sum_{d' \in D} q_d(d') V(T(\tilde{s}, d, a), d')\}, \end{aligned} \quad (5)$$

$\tilde{s} \in \tilde{S}, d \in D.$

#### IV. THE AGENT

In this paper we evaluate the performance of three agents based on different algorithms: a heuristic based on the  $(\zeta, Q)$ -Policy that we use as the baseline for performance evaluation, an approximate SARSA algorithm and an implementation of the REINFORCE algorithm.

**The  $(\zeta, Q)$ -Policy** based agent is not smart in the way that it does not learn over time. Due to the popularity of the  $(\zeta, Q)$ -Policy in practice we use it as a baseline for performance evaluation of the other agents. In the heuristic we iterate over  $s_1, \dots, s_k$  and replenish the respective warehouses by some amount  $a_i = Q_i$  if the current stock is below a level  $\zeta_i$  and there is still stock left in the factory  $s_0$ . At the end we set the production level for the next period to  $Q_0$  if  $s_0 - \sum_{i=1}^k a_i < s_0$

and 0 otherwise. The thresholds  $\varsigma$  and replenishment levels  $Q$  need to be set by the user when initializing the agent.

**Approximate SARSA** uses a linear approximation  $Q_\theta(\tilde{s}, a) = \theta^T \phi(\tilde{s}, a)$  of the Q-function. We chose this method as it solves the problem of exponentially growing state- and action- spaces that is likely to occur in multidimensional environments. Furthermore it allows us to use our knowledge of the environment (especially the structure of the reward function  $r(\tilde{s}, d, a)$ ) to design  $\phi(\tilde{s}, a)$  such that it preserves some of the MDPs structure. The algorithm is implemented as depicted in [3] where we use an  $\epsilon$ -Greedy strategy to select the action and a search-then-converge stepsize rule to update the learning rate  $\alpha$  and the probability of choosing a random action  $\epsilon$  as defined in (3).

One of the crucial tasks when designing the approximate SARSA agent is the model for  $\phi : \tilde{S} \times A \rightarrow \mathbb{R}^p, p \in \mathbb{N}$ . This can be understood similar to feature engineering where we use some state  $s$  and some action  $a$  to compute  $p$  features that should be capable of describing  $Q(s, a)$ . In our case we create over 15 different types of features that can be binary as well as continuous. One of the main ideas is to get a rough estimate of the next demand and state by

$$\begin{aligned} \hat{d}(\tilde{s}) &= d_{last} + (d_{last} - d_{hist}) = 2d_{last} - d_{hist}, \\ \hat{s}(\tilde{s}, a) &= T(\tilde{s}, \hat{d}(\tilde{s}), a)_{(d, d_{last})}, \end{aligned} \quad (6)$$

where  $T(\cdot)_{(d, d_{last})}$  is the transition function excluding its last two elements. This way the agent can get a basic understanding of rewards and penalties associated with  $\hat{s}$ . Among others we use the expected penalty costs, the expected reward and the squared difference of the stock levels  $s$  from their 25%, 50% and 75% maximum capacity. Furthermore we include the respective rewards and costs for two scenarios where  $\hat{d}^+ = \hat{d}(\tilde{s}) + 1$  and  $\hat{d}^- = \hat{d}(\tilde{s}) - 1$ . While designing the agent we ended up with two different functions  $\phi$  and  $\phi'$  that affect the agents behavior. Both functions have been tested and are included in our code. We encourage the reader to try out both. A full list of all features can be found in the appendix.

When testing the approximate SARSA algorithm we found that for some environments the parameters  $\theta$  would increase until computations became numerically unstable. To avoid this issue we restrict the temporal difference

$$\delta_n = r(\tilde{s}_n, d_n, a_n) + \gamma \theta_n^T \phi(\tilde{s}_{n+1}, a_{n+1}) - \theta_n^T \phi(\tilde{s}_n, a_n) \quad (7)$$

that is used to update  $\theta_n$  within the interval  $[-1e100, +1e100]$ . Furthermore we initialize the agent with a very small stepsize  $\alpha$ .

**REINFORCE** is based on a parametrized policy for which the expected reward has to be maximized. We assume our policy to be parametrized as a softmax function for multiple actions. The softmax function assigns a probability for each action and thus behaves as a stochastic policy which can, however, converge to a deterministic one. We discretize the action space such that we only have three available actions per entity (for each warehouse and the factory). For a warehouse, for example, these possible actions are an order of 0,1 or 2 full trucks from the factory. For a given state, some actions might not be available. Since a projection of an unfeasible

action to the closest feasible action is not possible in our environment, we used an alternative approach which only updates parameters for available actions. This leads to the following policy parametrization:

$$\begin{aligned} p(a = a^{(k)} | s) &= \pi_\Theta(a = a^{(k)}, s) = \\ &= \frac{e^{\Phi(s)^T w_{a^{(k)}}} \cdot f(a^{(k)} | s)}{\sum_{j=1}^{n_a} e^{\Phi(s)^T w_{a^{(j)}}} \cdot f(a^{(j)} | s)} := \sigma_k(s) \end{aligned} \quad (8)$$

where

$$f(a^{(j)} | s) = \begin{cases} 1 & \text{if } a^{(j)} \text{ is allowed in state } s \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

As we allow 3 different actions per entity, the number of possible actions is  $n_a = 3^{(n_{stores}+1)}$ . The feature map  $\Phi(s) : \mathbb{Z}^{n_{stores}+1} \rightarrow \mathbb{R}^d$  maps the state into a space of dimension  $d$  in which an optimal action is easily separable from the other ones. We tried several different options for  $\Phi(s)$  including first and second order terms and radial basis functions. The parameter matrix  $\Theta$  is assembled in the following way:

$$\Theta = (w_1 | w_2 | \dots | w_{n_a}) \in \mathbb{R}^{d \times n_a} \quad (10)$$

where  $w_i \in \mathbb{R}^d$  (including a bias). We initialized  $\Theta$  to zeros in order to start with equal probabilities for each action.

A gradient ascent method is used to find the parameters that maximize the expected reward following this policy. The gradient for our parametrization evaluates to:

$$\nabla_{w_i} \ln(\pi_\Theta(a^{(k)}, s)) D_t = \begin{cases} (1 - \sigma_i(s)) \Phi(s) D_t & \text{if } i = k \\ -\sigma_i(s) \Phi(s) D_t & \text{if } i \neq k \end{cases}$$

The full derivation for this gradient is shown in the Appendix. A key parameter for this algorithm is the learning rate which we left constant, but close to zero since a small learning rate allows a lot of exploration which is needed for all parameters to adapt.

## V. RESULTS AND DISCUSSION

### A. Performance of your Agent in your Environment

To test the algorithms, two main test were done. All of the tests were computed using 10000 episodes of 24 steps each in which our demand distributions make a whole cycle, allowing it to have the full seasonal demand.

The first test is a simple test, having only the factory and one warehouse, with no storage cost in the factory and small cost of production, transportation and storage in the warehouse. The objective of this test was to see if the agents can learn that even if producing, sending and storing is a bad decision in short term, because there is a penalty cost that accumulates for not sending the goods to the warehouse, the overall strategy should be to do it anyway.

But because there is a cost of production, the optimal strategy would be to not build as much as we can but to build enough to be able to satisfy the season of high demand but no more than that.

The graph seen in Figure 3 shows that both the approximate SARSA algorithm and the REINFORCE can learn this task,

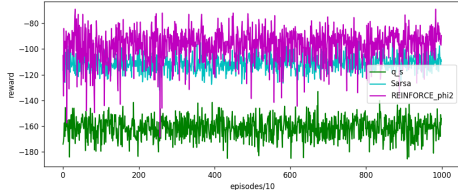


Fig. 3. Rewards over time for the three agents in the first test.

with the REINFORCE algorithm having slightly better results and with notable better results that the  $(\varsigma, Q)$ -Policy based agent, who sends and produces goods before they are needed and ends up getting a higher penalty. This shows that the basis of our environment can be learned by our agents even if they are only getting weak signals.

The second test is a harder environment. In this environment the number of warehouses went up to 3, each with different properties. The first warehouse has a high storage and transportation cost, the second warehouse has no storage cost but high transportation cost and the third one no transportation nor storage cost. Furthermore, because of how the production and the demand are balanced there are not enough goods to satisfy all the demand, the agent need to learn everything he learned in the simple scenario but also to prioritize some warehouses before others.

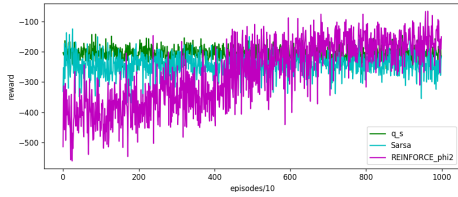


Fig. 4. Rewards over time for the three agents in the second test.

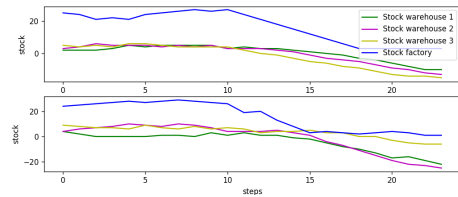


Fig. 5. Stocks for the  $(\varsigma, Q)$ -Policy based agent and the REINFORCE agent in the second test.

This second graphs seen in Figure 4 shows the  $(\varsigma, Q)$ -Policy based agent who begins with better results thanks to his simple solutions that allows him to produce a lot and send it to fill the stores one by one, but after all the training the REINFORCE algorithm manages to learn a better strategy. As shown in Figure 4, the storage levels in the different stores for the  $(\varsigma, Q)$ -Policy strategy are pretty similar while on the REINFORCE we have the warehouse 1 who has at beginning almost zero stock due to its high storage cost, and when the demand is too high both warehouse 1 and 2 are forfeit to the store three, thanks to his zero cost of storage

and transportation.

The SARSA on the other hand does not manage to learn the right strategy, building less stock during the low demand and not prioritizing early the warehouse three. The main problem from the not learning of the SARSA is the difficulty to choose the right features and to allow the agent to learn what needs to be learn with this features. Here, we cant just have a feature that is storage level due to the fact that having really high storage is bad but low is also bad. And because of how the value of each action in every state is calculated is a weight times the feature, he can either say producing is bad or producing is good. So a balance between the features needs to be reached, where the agent can say that producing is good unless there is too much. This balancing of feature is hard to learn and provoques bad results in too complicated scenarios. This might be solved with a different combination of features but for this particulars parameters none was found.

The readers are encouraged to create their own scenarios in testing.py to see the different performances of the agents. To compare their performances or recreate the graphs seen in this paper graph\_creation.py needs to be run.

### B. Performance of your Agent in the ALife Environment

Our implementation of the approximate SARSA algorithm uses a very specific value function approximation that depends strongly on assumptions about the underlying environment. Since this environment has no structural similarity with the ALife<sup>3</sup> environment it can not be used for testing. To adapt the algorithm to the ALife environment it would be necessary to create new features and design a different function  $\phi$  to be used within the existing agent.

In general, the REINFORCE approach is very nicely applicable to the Alife environment as was shown in the lab. Allowing multiple possible actions can also be a reasonable extension. Our implementation, however, is not directly transferable, since we use additional information about which actions are possible which is not available in the Alife environment.

## VI. CONCLUSION AND FUTURE WORK

The REINFORCE approach is certainly a viable option for this environment as it performs well in each scenario. Of course, the algorithm does have some downsides. As a gradient ascend algorithm it can easily get stuck in local optima which can easily arise in this possibly highly non-convex setting. Additional drawbacks come from the softmax parametrization which requires parameters for each action. This requires a lot of training and can lead to difficulties for actions that are rarely feasible. Moreover, the actions to replenish 0, 1 or 2 trucks are treated as completely separate, dropping the possibility to make use of their ordering. Another parametrization that would make use of this order would be e.g. a Gaussian setting. Unfortunately, our implementation of using independent Gaussians to parametrize the actions for each entity individually failed to show reasonable learning and was dropped therefore for this report.

<sup>3</sup><https://github.com/jmread/alife>

This section summarizes the paper: Your environment and agent, its strength and its weaknesses. Also remark about what would be the next steps you would take if you or someone else were to continue/extend this project. Note that for the initial submission you are limited strictly to 4 pages (double column), *not including references*. An extra page will be allowed for final submission (after the initial reviews).

## VII. APPENDIX

Derivation of the gradient for REINFORCE:

$$\nabla_{w_i} J(\Theta) = \nabla_{w_i} \ln(\pi_{\Theta}(a(t) = a^{(k)}, s(t))) D_t$$

where

$$\begin{aligned} & \nabla_{w_i} \ln(\pi_{\Theta}(a(t) = a^{(k)}, s)) = \\ & = \nabla_{w_i} (\Phi(s)^T w_{a^{(k)}}) + \nabla_{w_i} \ln(f(a^{(k)}|s)) + \\ & \quad - \ln\left(\sum_{j=1}^{n_a} e^{\Phi(s)^T w_{a^{(j)}}} \cdot f(a^{(j)}|s)\right) = \\ & = \Phi(s) + 0 - \frac{e^{\Phi(s)^T w_{a^{(k)}}} \cdot f(a^{(k)}|s)}{\sum_{j=1}^{n_a} e^{\Phi(s)^T w_{a^{(j)}}} \cdot f(a^{(j)}|s)} \cdot \Phi(s) = \\ & \quad = (1 - \sigma_i(s)) \cdot \phi(s) \end{aligned} \tag{11}$$

## REFERENCES

- [1] L. N. Moritz. Target Value Criterion in Markov Decision Processes, *Karlsruhe, KIT, Diss.*, 2014.
- [2] W. B. Powell, Approximate Dynamic Programming: Solving the curses of dimensionality, *John Wiley & Sons*, 2007.
- [3] J. Read. Lecture IX - Deep reinforcement learning. *INF581 Advanced Topics in Artificial Intelligence*, 2018.
- [4] K. Furmans, D. Arnold, Materialfluss in Logistiksystemen, *Springer Berlin Heidelberg*, 2009.
- [5] H. Tempelmeier, Inventory management in supply networks: problems, models, solutions, *Books on Demand*, 2011.