



## Manufacturing & Service Operations Management

Publication details, including instructions for authors and subscription information:  
<http://pubsonline.informs.org>

### Inventory Service-Level Agreements as Coordination Mechanisms: The Effect of Review Periods

Elena Katok, Douglas Thomas, Andrew Davis,

To cite this article:

Elena Katok, Douglas Thomas, Andrew Davis, (2008) Inventory Service-Level Agreements as Coordination Mechanisms: The Effect of Review Periods. *Manufacturing & Service Operations Management* 10(4):609-624. <https://doi.org/10.1287/msom.1070.0188>

Full terms and conditions of use: <https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact [permissions@informs.org](mailto:permissions@informs.org).

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2008, INFORMS

Please scroll down for article—it is on subsequent pages



With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

# Inventory Service-Level Agreements as Coordination Mechanisms: The Effect of Review Periods

Elena Katok, Douglas Thomas, Andrew Davis

Supply Chain and Information Systems Department, Smeal College of Business, Penn State University,  
University Park, Pennsylvania 16802 {ekatok@psu.edu, dthomas@psu.edu, amd361@psu.edu}

A supplier stocking goods for delivery to a retailer may face a (finite-horizon) service-level agreement (SLA). In this context, the SLA is a commitment by a supplier to achieve a minimum fill rate over a specified time horizon. This kind of SLA is an important, but understudied coordination mechanism. We focus on the impact of two contract parameters: the length of the review period and the magnitude of the bonus for meeting or exceeding the service-level target. For a supplier following a base stock (order-up-to) inventory policy, increasing the bonus increases optimal supplier stocking levels, whereas lengthening the review period may increase or decrease optimal stocking levels. We investigate these mechanisms in a controlled laboratory setting and find that longer review periods are generally more effective than shorter review periods in inducing higher stocking levels. As in several earlier laboratory studies, the explanation lies in the improved feedback reliability that longer review periods provide. The primary managerial implication of our findings is that, in practice, longer review periods may be more effective than shorter ones at inducing service improvements.

*Key words:* service-level agreements; behavioral operations management; supply chain management; experimental economics

*History:* Received: June 30, 2006; accepted: June 1, 2007. Published online in *Articles in Advance* January 4, 2008.

## 1. Introduction and Motivation

A supplier stocking goods for eventual delivery to a customer must trade off the negative consequences of stocking insufficient inventory with those of excess inventory. Consider the case in which the supplier charges a wholesale price to a retailer who then charges a retail price to the marketplace. The markup charged by the supplier causes the retailer to order less than the supply-chain-optimal quantity. Spengler (1950) was the first to note this *double-marginalization* problem. A variety of coordination agreements have been implemented or proposed in the literature to address this problem. Recently, Cachon (2003) reviewed several contractual mechanisms addressing this problem.

In this paper, we investigate (finite-horizon) service-level agreements (SLAs) as coordination mechanisms. These agreements are used to improve coordination by inducing suppliers to place higher orders. In an SLA, the supplier agrees to meet some predefined service level (typically the fraction of orders filled) over a specified review period. In some cases, there are

contractual financial penalties and rewards associated with failing or achieving a target service level for a particular time period. Another possibility is that the service level is part of a supplier scorecard used to evaluate supplier performance. In such a case, the negative consequences of failure may be more difficult to quantify. A recent Aberdeen survey (Kay 2005) reported that 70% of manufacturing companies declared that supplier performance, particularly on-time delivery and fill rates, is critical to their business operations.

An earlier analytical study has shown that finite-horizon SLAs can have negative consequences for both the supplier and customer in terms of long-run cost and profit if suppliers react optimally to SLAs (Thomas 2005). Furthermore, in terms of the agreement, the size of the penalty or bonus and the length of the review period can strongly affect the supplier's stocking decisions and, thus, customer performance, often in counterintuitive ways. In this paper, we consider the effects of these two variables on suppliers' stocking levels using controlled laboratory experi-

ments with human subjects. Our goal is to test the effects of the size of the bonus and the length of the review period on actual decisions, and to compare those decisions to theoretical (optimal) benchmarks.

We test the effects of the size of the bonus and the length of the review period in the laboratory because it allows us to induce assumptions made in the theory that ordinarily cannot be controlled in the field. Thus, when we observe differences between the actual behavior and theoretical predictions, we can attribute these differences to behavioral factors. There is a long tradition of using laboratory experiments in the decision-making literature. We refer the reader to Kagel and Roth (1995) for a review of various problems in economics that have been studied using experimental methods, and to Camerer (2003) for a review of the literature emphasizing the link between economics and psychology.

The use of laboratory experiments to study problems in operations management also has a long history (see Bendoly et al. 2006 for a review). Specifically relevant to our study is the work of Rapoport (1966, 1967), who found that decision makers in a stochastic multistage inventory task generally undercontrol the system, and although demand draws are independent, orders are correlated with past demand. More recently, Schweitzer and Cachon (2000) also found that in an even simpler, single-period inventory ordering task (the newsvendor problem), decision-makers tend to place orders that are correlated with past demand draws, even though demand draws are independent. Resulting average orders tend to be biased relative to the optimal orders in the direction of the average demand, and Schweitzer and Cachon (2000) note that this phenomenon is consistent with the “anchoring and insufficient adjustment” (p. 404) bias (participants start the game by anchoring on the average demand and then insufficiently adjust toward the optimal order). Anchoring and insufficient adjustment are consistent with two well-documented findings in the behavioral decision literature from the early days of the field: People have limited information processing capacity, and people are adaptive (e.g., Hogarth 1987 and references therein). Ben-Zion et al. (2007) replicate the Schweitzer and Cachon (2000) result for different demand distributions. Bostian et al. (2007) use an adaptive model

with reinforcement learning to explain this ordering behavior. Lurie and Swaminathan (2007) report that more frequent feedback sometimes actually degrades performance and slows down learning. Bolton and Katok (2007) find that the anchoring and insufficient adjustment bias persists with extended experience and under a variety of informational manipulations, but having decision makers place standing orders for multiple periods eliminates this bias. A bias similar to the anchoring and insufficient adjustment bias in the newsvendor problem has also been observed in some market entry and political participation experiments, and Goeree and Holt (2005) use bounded rationality to explain this behavior. Keser and Paleologo (2004) have shown that the anchoring on mean demand behavior is replaced by another anchoring mechanism, the tendency to split profit from sold units equally, in the simple newsvendor environment as soon as the supplier is included as a decision maker in the experimental game.

In the next section, we summarize the theoretical predictions of fill-rate performance over different-length review horizons with varying bonuses. We then describe our experimental design, which manipulates those parameters, and resulting research hypotheses (§3). We present our results and discuss how they relate to the research hypotheses in §4, and in §5 we summarize our results and discuss how they relate to both behavioral and analytical literature.

## 2. Analytical Results

### 2.1. General Model

For an inventory system with stationary demand and a stationary stocking policy, the long-run fill rate can be calculated by computing the expected units satisfied per period (or per replenishment cycle) and dividing this by the average demand. In a finite-horizon setting, the achieved fill rate is a random variable. Chen et al. (2003) and Banerjee and Paul (2005) investigate the behavior of the expectation of the achieved fill rate over a finite horizon; however, a supplier facing an SLA may be interested in the probability of meeting the specified target service level, rather than the expectation. Thomas (2005) investigates the distribution of the fill rate achieved over a finite horizon, including the probability of meeting a

specified target. It is worth noting that in all of those papers, as well as this one, the form of the inventory policy is restricted to a stationary, order-up-to policy. Such a policy is not necessarily optimal for a supplier facing a finite-horizon SLA; however, stationary policies are easy to implement and common in practice.

To focus on the implications of the SLA, we choose a simple, periodic-review inventory system with no ordering cost and zero lead time (next period delivery). Over a  $T$ -period horizon, the supplier faces demands  $D_i, i = 1, \dots, T$ . At the end of each period, the supplier incurs a holding cost  $h$  per unit held in inventory and a shortage cost  $p$  per unit for unfilled orders. In addition to those costs, the supplier receives a bonus if her fill rate over the  $T$ -period horizon meets or exceeds the threshold fill rate,  $\alpha_0$ . Because we will be making comparisons across different review horizon lengths, we will refer to the bonus amount in *per-period* terms. Let  $B$  denote the per-period bonus amount, implying a bonus of  $B \times T$  for the  $T$ -period horizon. To clarify, the supplier gets the entire bonus of  $B \times T$  if she achieves the target fill rate and zero otherwise.

Let  $S$  represent the supplier's order-up-to stocking level over the  $T$ -period horizon. The units of demand satisfied in any period  $t$  is then  $\min(D_t, S)$ . The supplier's cost function has two components. First, in each period there is the familiar expected holding and shortage costs

$$G_t(S) = pE(D_t - S)^+ + hE(S - D_t)^+.$$

We assume throughout our experiments that the demands are independent and identically distributed across periods; thus, we can drop the subscript  $t$  from  $G_t$  and represent the expected holding and shortage cost over the  $T$ -period horizon as  $\sum_t G(S) = TG(S)$ . Next, for a given  $S$ , the fill rate over the review horizon is

$$\alpha(S) = \frac{\sum_{t=1}^T \min(D_t, S)}{\sum_{t=1}^T D_t}.$$

The expected bonus to the supplier then is  $BT \Pr[\alpha(S) \geq \alpha_0]$ . We can now represent the suppliers expected cost *per-period* as

$$C(S) = G(S) - B \Pr[\alpha(S) \geq \alpha_0]. \quad (1)$$

The expected newsvendor cost (the first term in Equation (1)) is convex in the order-up-to level. The

expected bonus term is not necessarily convex, although it is monotonically nondecreasing. This monotonicity guarantees that any optimal solution to this problem has order-up-to level greater than or equal to the optimal solution to the newsvendor problem without the bonus, or more generally, the optimal order-up-to levels must be nondecreasing in  $B$ .

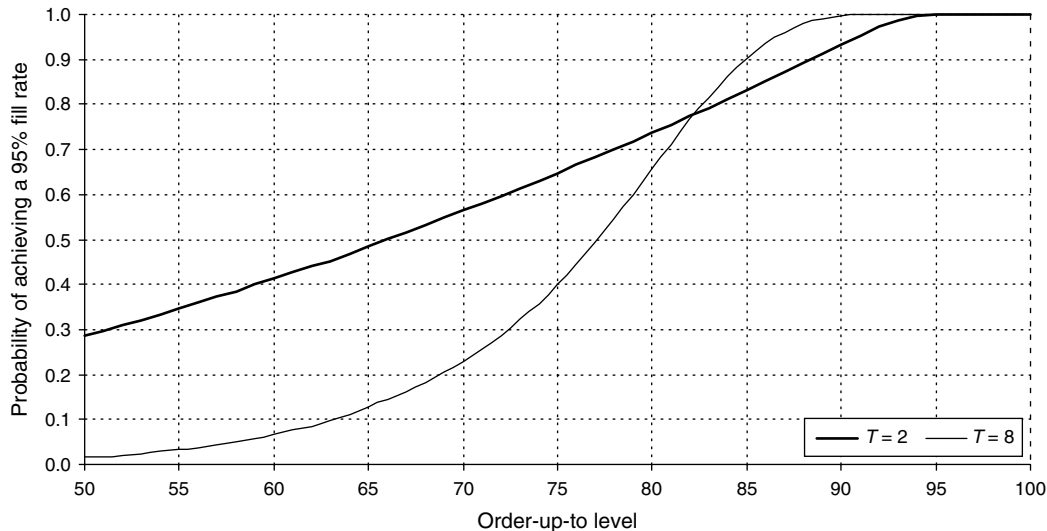
## 2.2. Laboratory Example

The discussion of the laboratory example is based on the analytical model in Thomas (2005). The standard expected newsvendor costs (per period) are unaffected by the review horizon length, whereas the expected bonus term does depend on the review horizon. Consider Figure 1, which shows the probability of meeting a 95% fill-rate target as a function of  $S$  for review horizons  $T = 2$  and  $T = 8$ , where per-period demand follows a discrete uniform distribution from 1 to 100. Note that at lower values of  $S$ , the probability of meeting the target with  $T = 2$  is substantially higher than for  $T = 8$ , whereas this is reversed for higher values. For this demand distribution,  $S = 78$  would result in a *long-run* fill rate of 95%, which means that as  $T$  becomes very large, the probability of meeting the 95% target goes from close to 0% to close to 100% at  $S = 78$ . We see this behavior starting to emerge in the  $T = 8$  curve. For long review horizons, the optimal value of  $S$  will either be close to the optimal stock level for the traditional newsvendor problem (when  $B$  is small), or close to the stock level where the "steep" part of the probability of meeting the target curve starts to flatten out (around 86 for the  $T = 8$  curve in Figure 1). When  $T$  is small, and the probability of meeting the target curve does not have dramatically "steep" and "flat" sections; the optimal value of  $S$  can take on many different values depending on  $B$ . We demonstrate this in the next section when we present our benchmark problems for the experiments.

## 3. Design of the Experiment

### 3.1. Methods

In our laboratory experiment, one human subject sets the order-up-to level for  $T$  periods at a time. Each *review period* consists of  $T$  *periods*. At the end of the  $T$  periods, the participants observe the actual demand

**Figure 1** Probability of Meeting 95% Fill Rate as a Function of Order-Up-To Level

during each period, the actual inventory and backlog levels and costs, the cumulative demand and inventory and backlog, the actual fill rate that resulted from the ordering policy, and the resulting profit. See the online appendix for sample instructions.<sup>1</sup>

In all of our treatments, the customer demand follows a discrete uniform distribution from 1 to 100 units per period, and both the holding and the backlog costs are set at 1 “franc” per unit.<sup>2</sup> We made the holding and backlog costs symmetric to create an environment in which, absent a bonus, the optimal  $S$  corresponds to the average demand—a fairly transparent solution.<sup>3</sup> The target fill rate is 95% in all treatments.

Our design manipulates two factors: We set the bonus levels at  $B = 0, 5, 25,$  or 50 francs per period, and we set the review period to  $T = 2$  or  $T = 8$ .

<sup>1</sup> An online appendix to this paper is available on the *Manufacturing & Service Operations Management* website (<http://msom.pubs.informs.org/ecompanion.html>).

<sup>2</sup> We use “franc” to refer to a unit of experimental currency. All explanations are presented to participants in terms of francs, and their earnings from the experiment are in francs. Francs are converted to U.S. dollars at a prespecified rate at the end of the experiment.

<sup>3</sup> This solution should be transparent to the reader, and it was our expectation that it would also be transparent to the participants in our experiment. Overall, deviations from this solution were not very large, but were significant (see the next section).

$T$  represents the number of inventory replenishments during the review horizon, and this value could range dramatically in practice. We choose small ( $T = 2$ ) and moderate ( $T = 8$ ) values for the experiments to facilitate investigation of the review horizon effect. In practice, the number of replenishment cycles during the review horizon may fall in this range (e.g., weekly replenishment with monthly performance reviews), but may also be somewhat larger (e.g., daily replenishments with quarterly reviews). Values of  $T$  substantially larger than those in our experiments would result in much less variability in fill-rate performance, making the experiment less informative.

In each treatment, subjects make 50 *ordering decisions* under one set of parameters, followed by another 50 decisions under a different set of parameters. This *within-subject design* has the advantage of increased statistical power because it automatically controls for individual differences across subjects (Camerer 2003, pp. 41–42). The main disadvantage of the within-subjects design is that, because participants have to complete two different tasks, it is important to test for the *order effects*. Order effects refer to the possibility that the experience in the first task might bias the behavior in the second task. The standard methods for checking for the order effects is to vary the order of the tasks for different subjects, and then compare the outcomes of a task for the participants who performed it first to the participants who performed it

**Table 1** Summary of Experimental Design

Bonus ( $B$ )	Time horizon ( $T$ )	
	$T = 2$ (short)	$T = 8$ (long)
$B = 0$ (no bonus)	Session 1: $N = 16$	
$B = 5$ (low bonus)	Session 2: $N = 20$	Session 3: $N = 20$
$B = 25$ (medium bonus)		
$B = 50$ (high bonus)	Session 4: $N = 20$	
$B = 50$ (Exec)	Session 5: $N = 8$	

*Notes.* Numbers inside table cells refer to sample sizes in each treatment. Exec refers to subjects with managerial experience.

second (Camerer 2003, p. 40). We used this method in all of the treatments and found no evidence of order effects.<sup>4</sup>

Table 1 summarizes our design. Each cell in the table represents a single session that includes the same group of participants completing the experimental task with two sets of parameters. We list sample sizes ( $N$ ) in each session in the table as well. For example, Session 1 included 16 participants who completed two games with zero bonus, one game in the  $T = 2$  condition and one in the  $T = 8$  condition. Session 2 included 20 (different) participants who also completed two games, both in the  $T = 2$  condition, one with  $B = 5$  and one with  $B = 25$ . Because participants in our study did not interact among themselves, each participant constitutes a single independent observation, which we will use as the main unit for our statistical analysis.

Participants in treatments corresponding to the first four rows of Table 1 were students, mostly undergraduates, from a variety of majors at Penn State University. Each individual participated in a single session only. They were recruited using the online recruitment system. Cash was the only incentive offered, and these participants were each paid a \$5 participation fee, plus an additional amount based on their performance. Average earnings for those participants, including the participation fee, were \$18. Each session lasted approximately 45 minutes. All

<sup>4</sup>Specifically, we checked for order effects by comparing the order levels between subgroups of subjects who used a specific set of parameters first versus second and found no statistical (or even meaningful) difference. All managers played the game in the  $T = 8$  condition first and the  $T = 2$  condition second, so no order-effect test was possible.

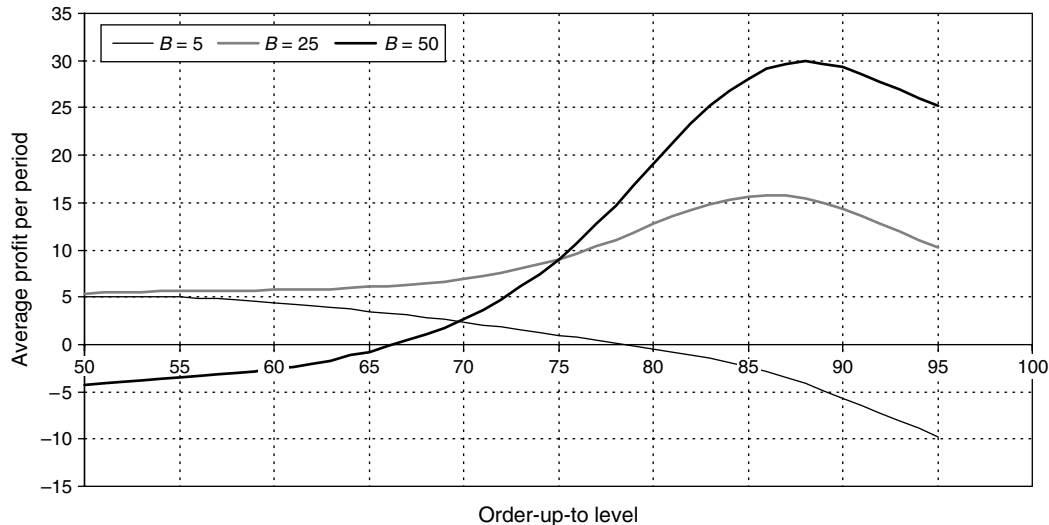
student sessions were conducted at the Laboratory of Economic Management and Auctions at Penn State University, Smeal College of Business, during April 2006.

Participants in the treatment corresponding to the last row in Table 1, labeled  $B = 50$  (Exec) were supply chain managers from firms affiliated with the Center for Supply Chain Research at Penn State University. These participants did not receive financial incentives. The session with professional managers was a replication of the otherwise identical session with student subjects. We conducted this session to check the effect of using subjects with managerial experience. Data for managers were collected during May and June of 2006. In total, 84 participants were included in our study: 76 students and 8 managers. All sessions were conducted using web-based software written using PHP and MySQL database back end.

To keep financial incentives constant across treatments, the bonus amount displayed in Table 1 represents the bonus per-period received in the event that the target fill rate was achieved during the review period; thus, in the two-period settings, the actual bonus amounts were 0, 10, 50, and 100, and in the eight-period settings, the actual bonus amounts were 0, 40, 200, and 400. To keep the game frame in the domain of gains rather than losses (Kahneman and Tversky 1979), each period in each treatment included an additional fixed endowment amount (meant to represent the net profit from selling the product). That is, the profit function our participants face is simply a constant minus the total cost Equation (1) presented in the previous section. We varied the endowment amount across treatments so as to keep actual earnings roughly the same in all sessions. Endowment amounts per period were 40 in  $B = 0$  (both  $T = 2$  and  $T = 8$  conditions), 40 in  $B = 5$  (both  $T = 2$  and  $T = 8$  conditions), 40 in  $B = 25$  ( $T = 2$  condition), 30 in  $B = 25$  ( $T = 8$  condition), and 20 in  $B = 50$  ( $T = 2$  and  $T = 8$  conditions).

### 3.2. Experimental Hypotheses

Our first two hypotheses are guided by the results from §2. We formulate these hypotheses as qualitative shifts of the order-up-to level due to (a) the bonus amount, and (b) the length of the review period. For each experimental treatment, we calculate expected profit as a function of the order-up-to

Figure 2 Expected Profit for  $T = 8$  Condition

levels. Calculating the expected profit for a given  $S$  requires evaluating the probability of meeting the target service level, which is a convolution of random variables. We use simulation to evaluate this convolution for all possible values of  $S$ .<sup>5</sup> In general, because expected profit is not necessarily unimodal, one cannot employ an efficient line-search algorithm for finding optimal stock levels. In our case, because we only have 100 possible decisions, we evaluate expected profit at all possible stock levels. Figures 2 and 3 show these expected profits, and we use them as theoretical benchmarks. We formulate an additional hypothesis about how we expect the actual behavior to deviate from theory, based on results from earlier studies about inventory ordering behavior.

The first hypothesis speaks to the effect the bonus amount should have on the average order-up-to level.

**HYPOTHESIS 1 (BONUS AMOUNT).** *Higher bonuses should cause higher average order-up-to levels for both review periods.*

Theoretical predictions for bonuses of 0, 5, 25, and 50 are 50, 54, 70, and 92 for  $T = 2$  and 50, 51, 86, and 88 for  $T = 8$ , respectively (see Figures 2 and 3).

<sup>5</sup> To evaluate the probability of meeting the target, we use Latin Hypercube sampling and a very large number of replications, so simulation errors are insignificant (see McKay et al. 1979).

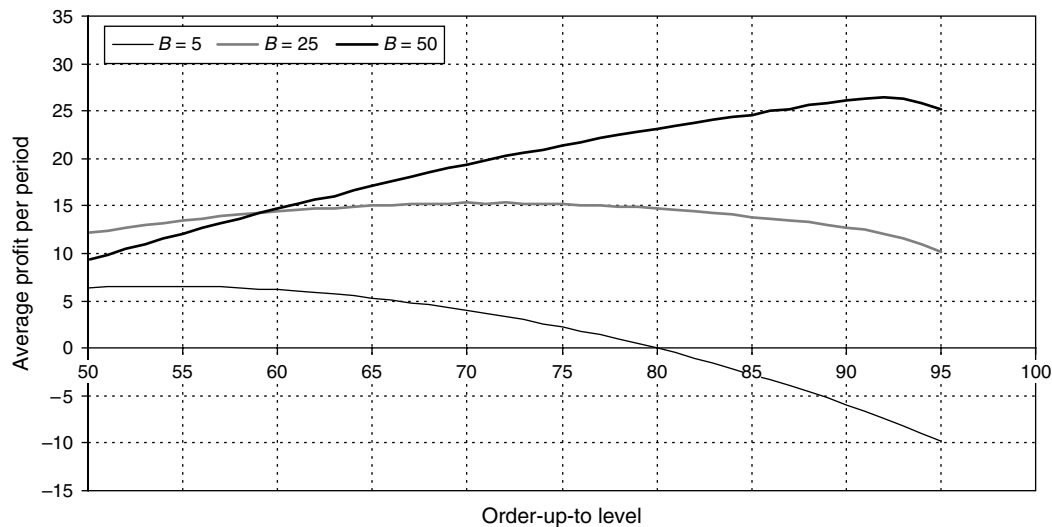
The second hypothesis speaks to the effect of the review period on the average order-up-to levels for various bonus amounts.

**HYPOTHESIS 2 (REVIEW PERIODS).** *The order-up-to levels depend on the length of the review period.*

When the bonus is 0, the review period has no effect on the order-up-to levels, which should be always 50; the bonus levels of 5 and 50 should induce lower average order-up-to levels when  $T = 8$  than when  $T = 2$  (decrease from 54 for  $T = 2$  to 51 for  $T = 8$  for  $B = 5$ , and from 92 for  $T = 2$  to 88 for  $T = 8$  for  $B = 50$ ); the bonus level of 25 should induce lower average order-up-to levels in the  $T = 2$  condition (70) than in the  $T = 8$  condition (86).

Whereas the first two hypotheses refer to theoretical benchmarks of how the order-up-to level should be affected by the bonus amount and by the length of the review period, the third hypothesis speaks to potential deviations between the actual behavior and theoretical predictions due to the anchoring and insufficient adjustment bias we mentioned in the introduction. In the context of our present game, absent a bonus, the optimal order and the average demand coincide. We intentionally selected these parameters to simplify the problem as much as possible, and to isolate the effect of the bonus on behavior. Generally, if decision makers anchor on average demand and then adjust (insufficiently) toward the

Figure 3 Expected Profit for  $T = 2$  Condition



optimal order, the implication for our game is that the actual order will be below optimal for all positive bonus amounts (with the caveat that, because in the  $B = 5$  condition the optimal order is very close to the average demand, differences may not be detectable).

**HYPOTHESIS 3 (ANCHORING AND INSUFFICIENT ADJUSTMENT).** *In conditions with positive bonuses, average orders should be below optimal (in these cases with optimal order quantity > the mean).*

## 4. Results

### 4.1. Descriptive Statistics

Table 2 summarizes the median, mean, and standard deviations of the order levels in each of our treatments, and compares them to theoretical benchmarks.

Because in the  $T = 2$  condition participants observed 100 periods and in the  $T = 8$  condition they observed 400, we also report descriptive statistics in the  $T = 8$  treatments for the first 12 decisions (corresponding to the first 96 of the 400 periods).<sup>6</sup> All the

<sup>6</sup>The median order-up-to levels for the first 12 decisions in the  $T = 8$  condition are very similar to the medians for all 50 decisions, and standard deviations are generally higher. Using the 12-decision data in the subsequent analysis makes no difference to any of the statistical comparisons we report, with one exception (see Result 6 below). The difference between the orders in the  $B = 50$  condition is significant when all 50 decisions are used, but is not significant when only the first 12 decisions are used.

comparisons we report in this section use median orders for individual subjects as the unit of analysis and all 50 decisions in the  $T = 8$  condition. All  $p$ -values we report below are two-sided. For one-sample tests comparing median orders to their theoretical benchmarks, reported in Table 2, we use the Wilcoxon signed-rank test. For two-sample tests, we use the Mann–Whitney test. Table 3 shows expected profits, standard deviation of profit, and probabilities of meeting the target level for optimal order levels, as well as average profit, and the fraction of time the target was met in the actual experiments.

**RESULT 1.** In the no-bonus condition, median orders are not statistically different from 50 and the review period length does not induce different median orders.

We cannot reject the null hypothesis that the median order-up-to level is 50. We also do not find any evidence that the median orders differ in the  $T = 2$  and  $T = 8$  conditions when  $B = 0$ .

Result 1 is not surprising given that the optimal solution is equal to the average demand when  $B = 0$ . We are interested in exploring how the inclusion of a bonus changes orders from this baseline.

**RESULT 2.** Higher bonuses induce higher average order-up-to levels.

We can reject all the null hypotheses that median orders for different bonus levels and time horizons are equal (all  $p$ -values < 0.01).



**Table 2** Summary of Average and Median Order-Up-To Levels and Their Standard Deviations in All Treatments, As Well As Corresponding Theoretical Benchmarks and Results of Hypothesis Tests

Bonus		Time horizon							
		2				8			
		Data (50 decisions)	Theory	<i>p</i> -value*		Data (50 decisions)	Data (12 decisions)	Theory	<i>p</i> -value*
0	Median	51.44	50	0.1354	50.91	50.58	50	0.1375	0.7061
	Mean	52.31			51.53	52.29			
	Std dev	(4.63)			(5.16)	(6.01)			
5	Median	56.00	54	<b>0.0217</b>	60.20	61.79	51	<b>0.0005</b>	0.4328
	Mean	58.55			62.78	63.87			
	Std dev	(8.04)			(11.65)	(13.92)			
25	Median	68.57	70	0.1968	79.66	76.75	86	<b>0.0002</b>	<b>0.0001</b>
	Mean	68.19			79.04	76.64			
	Std dev	(8.45)			(7.54)	(7.70)			
50	Median	81.92		<b>0.0003</b>	86.58	85.63		<b>0.0415</b>	<b>0.0425</b>
	Mean	81.72			86.03	83.47			
	Std dev	(8.43)			(4.70)	(6.94)			
50 (Exec)	Median	83.70	92	<b>0.0391</b>	82.24	76.71	88	<b>0.0039</b>	0.9314
	Mean	83.24			81.65	75.77			
	Std dev	(9.55)			(3.97)	(7.30)			

Note. *P*-values below 0.05 are bold.

\*Wilcoxon signed-rank test, null hypothesis ( $H_0$ ): Median order = theoretical prediction.

\*\*Mann–Whitney test,  $H_0$ : Median order for  $T = 2$  condition equals median order for  $T = 8$  condition.

The second result suggests that SLAs are at least directionally effective in addressing the double-marginalization problem.

The next four results address the effect of the review horizon length on order-up-to levels under different bonus amounts.

**RESULT 3.** In the  $B = 5$  condition, the median order-up-to levels are above theoretical predictions, and, contrary to theoretical predictions, they are not different for  $T = 2$  and  $T = 8$  conditions.

We can reject the null hypothesis that the median order-up-to level is not different from 54 for  $T = 2$  and 51 for  $T = 8$ . We do not find any evidence that the median orders differ in the  $T = 2$  and  $T = 8$  conditions.

**RESULT 4.** In the  $B = 25$  condition, the median order-up-to levels are not different from theoretical predictions in the  $T = 2$  condition, and are below theoretical predictions in the  $T = 8$  condition. Consistent with the theory, order-up-to levels in the  $T = 2$

**Table 3** Summary of Profit and Probability of Meeting Fill-Rate Target at Optimal and Actual Order-Up-To Levels ( $S$ )

		$T = 2$					$T = 8$				
		Bonus	Mean	Stdev	CV	$S$	Prob. of success (%)	Mean	Stdev	CV	$S$
At optimal $S$	0	30.00	20.60	0.69	50	28.63	120.01	40.83	0.34	50	1.60
	5	33.10	21.55	0.65	54	33.33	120.74	41.62	0.34	51	1.72
	25	50.56	31.74	0.63	70	56.42	125.92	87.12	0.69	86	93.49
	50	52.84	41.56	0.79	92	97.35	239.11	93.54	0.39	88	97.93
Actual	0	26.13	53.34	2.04	52.31	33.25	107.74	49.12	0.46	51.53	4.88
	5	28.13	26.93	0.96	58.55	43.50	101.34	76.81	0.76	62.78	25.00
	25	32.47	38.84	1.20	68.19	62.55	103.36	100.10	0.97	79.04	70.00
	50	47.83	45.77	0.96	81.72	79.50	199.46	145.27	0.73	86.03	88.10
	50 (Exec)	39.88	46.83	1.17	83.24	76.89	170.92	164.43	0.96	81.65	74.40

Note. CV, Coefficient of variation.

condition are below order-up-to levels in the  $T = 8$  condition.

We test whether the median order-up-to level is 70 for  $T = 2$  and 86 for  $T = 8$ , and find that we cannot reject the null hypothesis for  $T = 2$ , but can reject it for  $T = 8$ . Consistent with theoretical predictions, the median order in the  $T = 2$  condition is significantly below the median order in the  $T = 8$  condition.

RESULT 5. In the  $B = 50$  condition, the median order-up-to levels are below theoretical predictions in both  $T = 2$  and  $T = 8$  conditions. This is true with student participants as well as managers.

We test whether the median order-up-to level is 92 for  $T = 2$  and 88 for  $T = 8$ , and find that we can reject the null hypothesis in both cases.

These last three results indicate that subjects' order-up-to levels are above theoretical predictions for the zero and low-bonus ( $B = 5$ ) treatments, and at or below theoretical predictions for medium- and high-bonus treatments ( $B = 25, 50$ ). The last two observations are consistent with a form of insufficient adjustment. From Table 3 we see that the probability of meeting the target at the *optimal* order level for  $T = 2$ ,  $B = 50$  is 97.35%, although the subjects only obtained the bonus 79.50% of the time.

Figure 4 shows the distribution of achieved fill rate at median order levels for the high ( $B = 50$ ) bonus, and may offer some insight to the underordering behavior for  $T = 2$ . Note that at the median order level for  $T = 2$ , subjects have a 67% chance of meeting

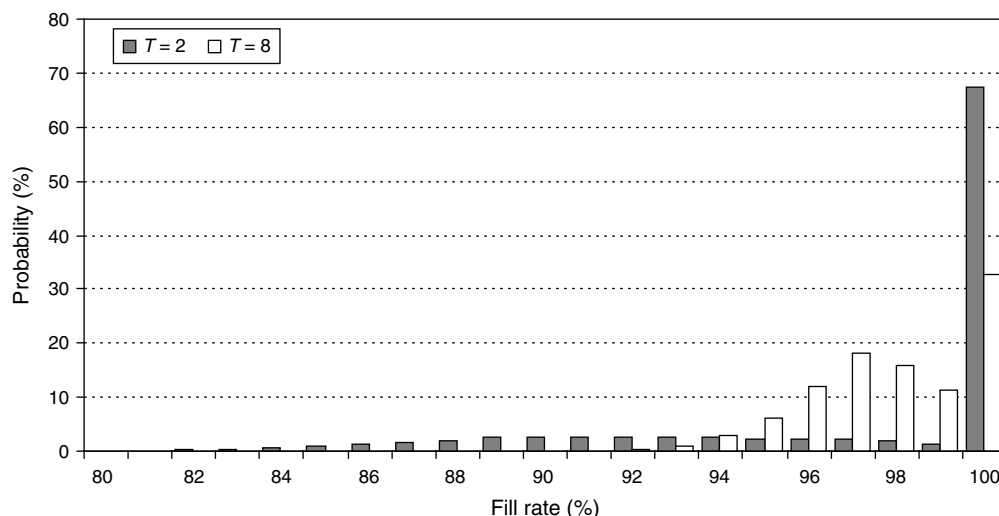
all demand (100% fill rate), and then some smaller chance of fill rates ranging from 82% to 99%. At this median level, the fact that a 100% fill rate is often observed may be discouraging subjects from appropriately increasing their stock level. For the  $T = 8$  case, the optimal and median order levels are 88 and 86, with probabilities of meeting the target fill rate of 98% and 88%, respectively (see Table 3). Examining Figure 4, we see that there is substantially less variability in the achieved fill rate for this longer review horizon, perhaps partially explaining why subjects do a better job of closing in on the optimal solution.

RESULT 6. In the  $B = 50$  treatment, order-up-to levels in the  $T = 2$  condition are below order-up-to levels in the  $T = 8$  condition with student participants, and are not significantly different with managers.

Contrary to theoretical predictions, the median order in the  $T = 2$  condition is significantly below (rather than above) the median order in the  $T = 8$  condition for students (this is the only result that does not hold when we use the first 12 instead of all 50 decisions). This observation is likely due to the more dramatic underordering that occurs for the  $B = 50$ ,  $T = 2$  treatment discussed above. There is no significant difference for managers. Note that the small sample size in the manager treatment may account for the lack of significance.

RESULT 7. Generally, managers' order-up-to levels are not closer to theoretical predictions than students'.

Figure 4 Distribution of Achieved Fill Rate at Median Order Levels in  $B = 50$  Treatment



We find no difference in median order-up-to levels for managers and students in the  $T = 2$  condition ( $p = 0.9062$ ), and the differences in the  $T = 8$  condition are significant ( $p = 0.0196$ ), but indicate that students' median orders are actually closer to theoretical benchmarks than managers'.

#### 4.2. Dynamic Behavior

We now examine ordering behavior over time to gain additional insights into how participants learn. We will start by looking at the order levels over time graphically, and then present a formal statistical analysis using a regression model. It is instructive to start by looking at the  $B = 0$  data because in those treatments, we would argue, the optimal solution is transparent. Figure 5 plots the average order-up-to level for the 50 decisions for the  $T = 2$  and  $T = 8$  treatments.

The first remarkable observation is that, in this (admittedly trivial) game, participants do not immediately recognize 50 as the optimal order. On the contrary, to the extent there is an “anchor,” it appears to be above 50, and it takes the participants in both treatments about 15 decisions to adjust their average order-up-to levels to be indistinguishable from 50. Even at that point, however, there is some variability in orders.

Figure 6 shows average orders over time in the  $B = 5$  treatments. Both sets of data appear to be above their respective optimal order levels. In the  $T = 2$  condition, the average order is 58.55, with a standard deviation of 8.04, and this average is not statistically above the optimal order of 54. In the  $T = 8$  condition,

Figure 5 Average Orders Over Time in  $B = 0$  Treatments

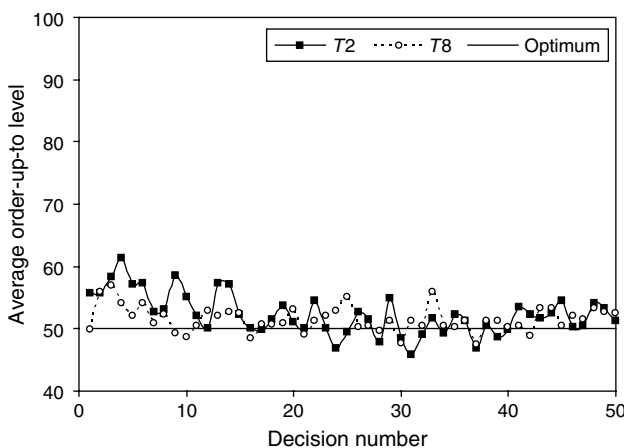
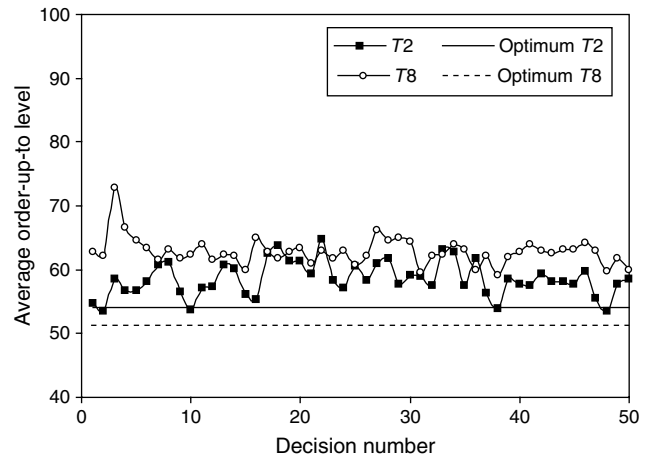


Figure 6 Average Orders Over Time in  $B = 5$  Treatments

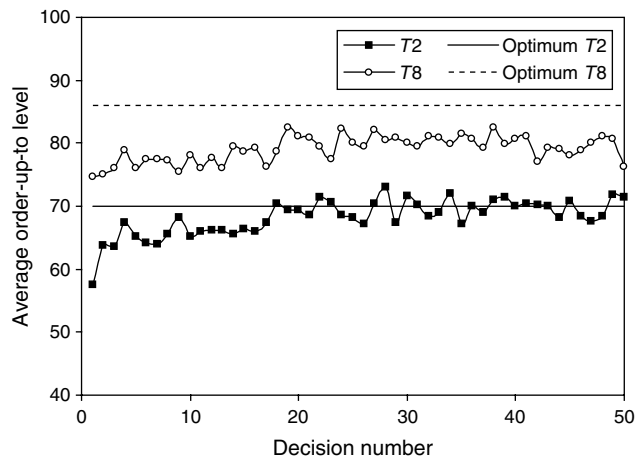


the average order is 62.78, with a standard deviation of 11.65, and this average is statistically above the optimal order of 51.

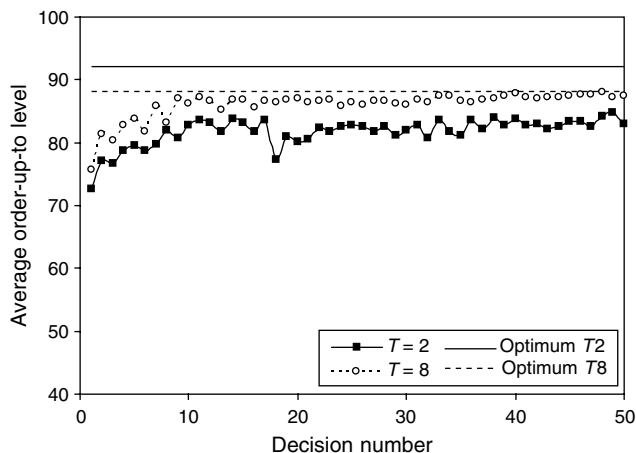
Again, as in the  $B = 0$  treatments, the average orders deviate from optimal orders in the direction away from the average demand. Figure 7 shows the same data for the  $B = 25$  treatments, and here we observe orders that start low and then increase in the direction of optimality. In the  $T = 2$  condition, the average orders converge to the optimal level, but in the  $T = 8$  condition they remain significantly below.

Figures 8 and 9 show the data for the  $B = 50$  treatments for student subjects and managers. In both of these figures we observe a similar pattern of orders starting out low and then increasing in the direction

Figure 7 Average Orders Over Time in  $B = 25$  Treatments



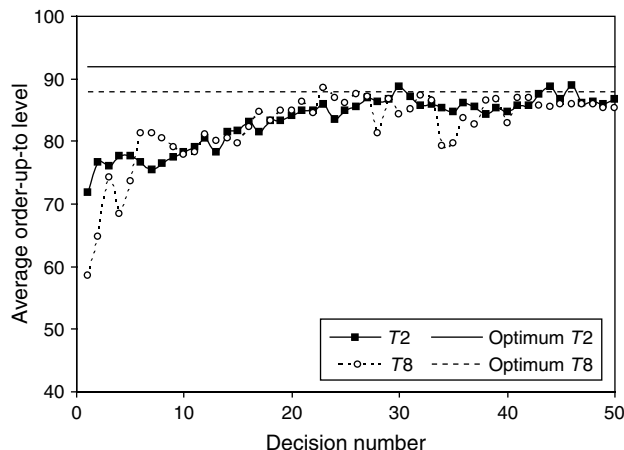
**Figure 8** Average Orders Over Time in  $B = 50$  Treatments; Student Subjects



of the optimal order-up-to level. However, in all treatments, average orders level off before they reach the optimal (profit-maximizing) levels. Whereas we see clear separation between the  $T = 2$  and the  $T = 8$  data in the student session (the average orders in the  $T = 8$  condition appear to be consistently higher), we see no such clear separation in the managers' data. Generally, although exhibiting the same general pattern as the students' data, managers' data are more variable and the learning appears slower. This may well be due to the small sample size.

The graphs in Figures 6–9 indicate some patterns of how ordering behavior evolves over time and is affected by the bonus levels and review-period

**Figure 9** Average Orders Over Time in  $B = 50$  Treatments; Manager Subjects



length. To measure these effects more systematically and compare their magnitudes, we use the following regression model

$$\begin{aligned} ORDER_{it} = & \beta_t \times t + \beta_{(t \times H)} \times (t \times HIGHBONUS) \\ & + \beta_B \times B + \beta_D \times DEMAND_{t-1} \\ & + \beta_{D8} \times (DEMAND_{t-1} \times (T = 8)) \\ & + \beta_T \times T + \beta_M \times MGR + \mu_i. \end{aligned} \quad (2)$$

$ORDER_{it}$  is the order-up-to level of participant  $i$  in decision number  $t$ , and Table 4 describes the explanatory variables in Model (2), along with the model coefficients we estimated (the  $\beta$ s) using the ordinary least-squares regression with fixed effects for individuals.

The estimates of Model (2) confirm several of the results we stated earlier. For example, the coefficient on  $B$  is positive and significant, meaning that higher bonus amounts do induce higher orders (Result 1). Also, the positive and significant coefficients on  $T$  confirm that the orders in the  $T = 8$  condition are higher than the orders in the  $T = 2$  condition (Result 6). The coefficient of the  $MGR$  variable is not significantly different from 0, confirming that there is no detectable effect due to subject pool with managerial experience (Result 7).

Model (2) allows us to formulate several additional results.

**RESULT 8.** Orders increase over time when the bonus is sufficiently high.

In low-bonus conditions ( $B = 0$  and  $5$ ) the overall time trend is negative, as can be seen by the coefficient on  $t$  being negative and significant. However, this trend reverses in the high-bonus conditions ( $B = 25$  and  $50$ ), as can be seen by the coefficient on  $t \times HIGHBONUS$ , which is positive, significant, and almost three times as large as the coefficient on  $t$ .

**RESULT 9.** Orders are correlated with last period's demand in the  $T = 2$  condition, but this demand-chasing behavior is virtually eliminated in the  $T = 8$  condition.

We draw this conclusion based on the coefficient for  $DEMAND_{t-1}$ , which is positive and significant, and the coefficient on  $DEMAND_{t-1} \times (T = 8)$ , which is negative, significant, and is of approximately the same magnitude as the coefficient on  $DEMAND_{t-1}$ .

**Table 4** Description of the Explanatory Variables and Model Estimates

Variable	Description	Coefficient (standard error)
$t$	Decision number 1 to 50	-0.11* (0.0112)
$t \times HIGHBONUS$	Interaction effect between decision number and the indicator variable, which is 1 for conditions $B = 25$ and $B = 50$ conditions and 0 for the $B = 0$ and $B = 5$ conditions.	0.28* (0.0138)
$B$	Bonus amount per period (0, 5, 25, or 50)	0.19* (0.0201)
$DEMAND_{t-1}$	Cumulative demand during the previous review period	0.05* (0.0039)
$DEMAND_{t-1} \times (T = 8)$	Interaction variable between last decision's cumulative demand and the $T = 8$ condition.	-0.04* (0.0043)
$T$	Review-period length	0.40* (0.1541)
$MGR$	Indicator variable, which is 1 for participants who are executives	1.63 (2.4160)
$\mu_i$	Individual participant fixed effect (average)	54.52* (1.3949)
$\hat{R}^2$	Adjusted $R$ -squared	0.6488

\* $p < 0.01$ .

RESULT 10. The orders anchor at a point that is somewhat above average demand.

The mean of the fixed-effects coefficients is 54.52, and it is significantly higher than the average demand of 50 ( $p = 0.0017$ ).

### 4.3. Learning

In this section, we further examine the mechanics participants use to adjust their stock levels. We use the learning direction theory (Selten and Stoeker 1986) as the basis for our analysis. In the context of our setting, the learning direction theory makes a clear prediction about how the order levels should be adjusted over time based on the outcome of the last period's decision: Following a period in which the target service level was met, the participants should decrease their order levels, and following the period in which the target service level was not met, participants should increase their order levels. Table 5 summarizes the actual adjustment behavior in treatments with positive bonuses.

The Table 5 cells in bold indicate behavior consistent with the learning direction theory. In each treatment after not meeting the target, participants are more likely to increase their order levels than to decrease them (binomial test,  $p < 0.001$  in all cases). After meeting the target, participants are more likely to decrease their order level than to increase them (binomial test,  $p < 0.05$ ). If we compute the proportion of adjustments consistent with the learning direction theory for each individual participant, we also find that participants are more likely to increase their

order-up-to levels after not meeting the target and to decrease them after meeting the target. The only exception is treatment  $T = 8$ ,  $B = 50$ , in which we find no statistical difference.

In Figure 10, we plot the average adjustment size over time for all the treatments.

It is clear from the figure that the overall size of adjustment decreases over time—participants learn more in the beginning of the game than at the end. To capture this effect more formally, we fit the following regression model

$$\begin{aligned}
 & |ORDER_{i,t-1} - ORDER_{it}| \\
 &= \beta_t \times t + \beta_{T8} \times (T = 8) + \beta_{B5} \times (B = 5) \\
 &\quad + \beta_{B25} \times (B = 25) + \beta_{B50} \times (B = 50) \\
 &\quad + \beta_M \times MGR + \beta_{t,T8} \times [t \times (T = 8)] \\
 &\quad + \beta_{t,B5} \times [t \times (B = 5)] + \beta_{t,B25} \times [t \times (B = 25)] \\
 &\quad + \beta_{t,B50} \times [t \times (B = 50)] \\
 &\quad + \beta_{t,M} \times [t \times MGR] + \mu_i.
 \end{aligned} \tag{3}$$

The dependent variable  $|ORDER_{i,t-1} - ORDER_{it}|$  is the absolute change in order level from decision  $t - 1$  to decision  $t$  for participant  $i$ . The independent variables are the decision number  $t$ , indicator variables for the long time horizon ( $T = 8$ ), the different bonus amounts ( $B = 5, 25, 50$ ), the executive subject pool ( $MGR$ ), and the interaction effects between the decision number and the rest of the explanatory variables. We summarize the estimates of Model (2) in Table 6.

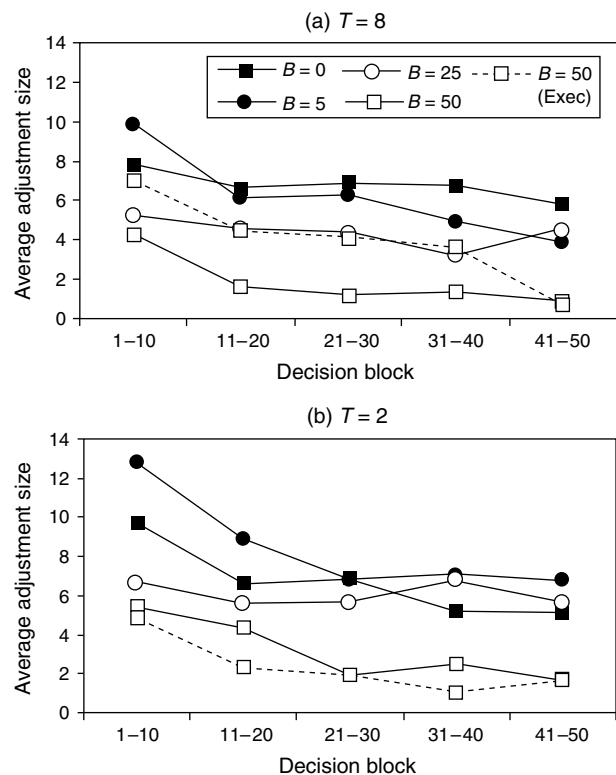
**Table 5** Number and Proportion of Adjustments in Response to Meeting or Not Meeting the Target Service Level

Treatment	Order adjustment		
	Decrease	Increase	Unchanged
<i>T</i> = 8; <i>B</i> = 5			
Did not meet target	112 (23.43%)	<b>163</b> ( <b>34.10%</b> )	203 (42.47%)
Meet target	<b>118</b> ( <b>38.56%</b> )	60 (19.61%)	128 (41.83%)
<i>T</i> = 2; <i>B</i> = 5			
Did not meet target	77 (20.16%)	<b>181</b> ( <b>47.38%</b> )	124 (32.46%)
Meet target	<b>192</b> ( <b>32.11%</b> )	131 (21.91%)	275 (45.99%)
<i>T</i> = 8; <i>B</i> = 25			
Did not meet target	99 (20.29%)	<b>214</b> ( <b>43.85%</b> )	175 (35.86%)
Meet target	<b>217</b> ( <b>44.11%</b> )	109 (22.15%)	166 (33.74%)
<i>T</i> = 2; <i>B</i> = 25			
Did not meet target	13 (7.60%)	<b>71</b> ( <b>41.52%</b> )	87 (50.88%)
Meet target	<b>185</b> ( <b>22.87%</b> )	112 (13.84%)	512 (63.29%)
<i>T</i> = 8; <i>B</i> = 50			
Did not meet target	176 (26.91%)	<b>204</b> ( <b>31.19%</b> )	274 (41.90%)
Meet target	<b>56</b> ( <b>43.08%</b> )	26 (20.00%)	48 (36.92%)
<i>T</i> = 2; <i>B</i> = 50			
Did not meet target	57 (18.75%)	<b>142</b> ( <b>46.71%</b> )	105 (34.54%)
Meet target	<b>222</b> ( <b>32.84%</b> )	150 (22.19%)	304 (44.97%)
<i>T</i> = 8; <i>B</i> = 50 (Exec)			
Did not meet target	123 (19.40%)	<b>216</b> ( <b>34.07%</b> )	295 (46.53%)
Meet target	<b>131</b> ( <b>37.86%</b> )	67 (19.36%)	148 (42.77%)
<i>T</i> = 2; <i>B</i> = 50 (Exec)			
Did not meet target	14 (9.52%)	<b>66</b> ( <b>44.90%</b> )	67 (45.58%)
Meet target	<b>188</b> ( <b>22.57%</b> )	151 (18.13%)	494 (59.30%)

As in Model (2), we estimate Model (3) using ordinary least squares with fixed effects for individuals.

The negative and significant coefficient on *t* indicates that the adjustments participants make to their orders decrease over time. Recall that demand-chasing behavior is prevalent in the short-review-period condition, but not in the long-review-period condition. Consistent with the demand-chasing observations, the

**Figure 10** Size of Adjustment Over Time



coefficient for *T* = 8 is negative and significant. We observe larger adjustments for the low-bonus (*B* = 5) condition and smaller adjustments for the medium-bonus (*B* = 25) condition.

#### 4.4. Summary of Results

We conclude the results section with a summary of our findings as they pertain to our three research hypotheses. We find strong support for Hypothesis 1. The higher bonus amounts do induce higher orders. We do not find general support of Hypothesis 2. The relationship between the order amount and the length of the review period does not generally confirm theoretical predictions. In fact, we find evidence that for medium and high bonuses, longer review periods induce higher orders, independent of the bonus. For low bonuses, orders under both review periods tend to be above their theoretical benchmarks and not different from one another. Of course, because our laboratory study included review periods of only two different lengths, the generality of this conclusion is limited.

**Table 6** Description of Explanatory Variables and Model Estimates for Model (3)

Variable	Description	Coefficient (standard error)
$t$	Decision number 2 to 50	-0.09** (0.0170)
$T = 8$	Indicator variable for treatments with the long time horizon ( $T = 8$ )	-1.25** (0.4586)
$B = 5$	Indicator variable for the $B = 5$ condition	5.47** (0.5675)
$B = 25$	Indicator variable for the $B = 25$ condition	-4.20** (0.8072)
$B = 50$	Indicator variable for the $B = 50$ condition	-0.90 (0.8782)
$MGR$	Indicator variable, which is 1 for participants who are executives	0.51 (0.8046)
$t \times (T = 8)$	Interaction variable between the decision number and $T = 8$	0.016 (0.0136)
$t \times (B = 5)$	Interaction variable between the decision number and $B = 5$	-0.09** (0.0209)
$t \times (B = 25)$	Interaction variable between the decision number and $B = 25$	0.04* (0.0205)
$t \times (B = 50)$	Interaction variable between the decision number and $B = 50$	-0.03 (0.0205)
$t \times MGR$	Interaction variable between the decision number and the executive subject pool	-0.07** (0.0225)
$\mu_i$	Average of individual fixed effects—the estimate of the average initial average adjustment amount.	10.58** (0.4613)
$\hat{R}^2$	Adjusted $R$ -squared	0.1909

\* $p < 0.10$ ; \*\* $p < 0.05$ .

For medium bonuses ( $B = 25$ ), the optimal order for  $T = 2$  is 70 and for  $T = 8$  it is 86, so the fact that the longer review period induces higher ordering levels is consistent with the theoretical prediction. However, for high bonuses ( $B = 50$ ), higher orders in the  $T = 8$  conditions are the opposite of the theoretical prediction. In those treatments, participants start by ordering too low and adjust in the direction of the optimal order, but because of the high probability of observing the fill rate of 100% in the  $T = 2$  treatment, they do not adjust their orders sufficiently. In the  $T = 8$  treatment, the feedback is more reliable and, consequently, participants are able to come closer to the optimal order.

One way to bring a common framework to our results is by noting a combination of behavioral effects. As we can see from Table 4, a bonus impacts players' orders in two ways: (1) in the initial order, and (2) in the direction of the order-up-to level adjustment. The anchoring effect on the initial order may be influenced by the availability heuristic (Tversky and Kahneman 1973) because a lump-sum bonus is more "available" than the holding or backorder costs, causing players to adjust orders accordingly. Because higher bonuses also increase optimal order-up-to levels, this has an impact on the location of players' orders relative to theoretical predictions. Players seem to overadjust the order-up-to level amount when given small bonuses and underadjust when given

higher bonuses. When no bonus is present, no adjustment takes place.

This brings us to Hypothesis 3, anchoring and insufficient adjustment, for which we do find some support. Interestingly, our data suggest that the anchor is slightly above average demand. The adjustment from the anchor is toward the optimal order (the adjustment is down in the low-bonus conditions and up in the high-bonus conditions), which is consistent with the hypothesis. Also consistent with the hypothesis, the adjustment is often insufficient; in many of the treatments, the average orders never reach the levels of the optimal order (see Figures 6–9).

In fact, our conclusions about the location of the anchor are not incompatible with earlier studies that documented the pattern of anchoring and insufficient adjustment. Schweitzer and Cachon (2000) and Bolton and Katok (2007) report that participants do somewhat better in solving the newsvendor problem when the optimal solution implies orders above average demand than when the optimal solution implies orders below average demand. Both of those studies suggested that this difference in performance may be caused by the perception that stocking out carries a negative connotation. Over time, participants respond to economic incentives and adjust their orders in the right direction.

We find evidence for demand chasing under  $T = 2$  but not under  $T = 8$ , which may have to do with

the perceived higher saliency of the  $T = 8$  decision. Decisions in the  $T = 2$  review periods and potential consequences of poor decisions “stick” for a relatively short amount of time (only two periods). Because bad consequences go away quickly and the probability of observing 100% fill rates is high (Figure 4), players can experiment with orders more leisurely, which could explain the “order-chasing” phenomenon here. Decisions under  $T = 8$ , however, “stick” for a longer time, and their consequences seem more salient. Because people consistently place too much weight on salient and tangible features of the environment (Kahneman et al. 1982, Taylor and Fiske 1975), players are less likely to experiment with orders, therefore avoiding demand chasing, and are more likely to be cautious, placing orders that try to meet the bonus, when it is present.

In summary, understanding the results in all of our treatments requires looking at both the effect of bonuses and the effect of demand chasing. Orders generally start at above average demand. Longer review periods decrease the variability of feedback, reduce demand chasing, and allow subjects to focus on the bonus. When the bonus is high enough, this causes orders to increase in the direction of the optimal order and performance improves. When the bonus is low, focusing on the bonus prevents orders from going down, so orders remain too high. When the bonus is 0, there is no focusing on the bonus and orders go down. This dynamic slows when the review period is short because demand chasing keeps subjects anchored closer to average demand and, also, the high probability of observing the 100% fill rate, even when the order was too low, causes orders not to increase when they should.

## 5. Discussion and Conclusions

Inventory SLAs are often used in practice to coordinate supply chains, and our data confirm that these mechanisms are directionally effective. High bonus amounts that reward suppliers for maintaining target fill rates indeed induce higher orders. But an important and little-studied parameter in these types of coordinating mechanisms is the length of the review period.

In theory, Thomas (2005) showed that the length of the review period matters, and the relationship between the optimal order-up-to level and the length

of the review period is complicated: If suppliers maximize expected profit, there are situations where shorter review periods would be preferred. Our laboratory experiment was designed to test these theoretical predictions, but we found that, in all of the settings in our study, including the setting with a very high bonus, a longer review period induced higher orders. We replicated the session with a very high bonus with executives and found that, qualitatively, their decisions were quite close to those of our student participants.<sup>7</sup>

Why do longer review periods induce higher orders, even when the optimal order is higher for shorter review periods? An explanation has to do, in part, with the difference in reliability of feedback that decision makers receive and the nature of this feedback. Bolton and Katok (2007) found that limiting decision makers to placing standing orders in the newsvendor problem dramatically improves performance by reducing the “demand-chasing” behavior, and we also observe this reduction in demand chasing in our study. Bolton and Katok (2007) note that this improvement is due to reducing the variability of feedback and allowing participants to “learn by doing.” In our setting, longer review periods also reduce the variability of feedback, and for the same reason as in the Bolton and Katok (2007) study, this improvement in reliability of information reduces demand chasing and allows participants to focus on the bonus. In treatments with a medium and high bonus, when initial orders are too low, focusing on the bonus causes orders to increase, which translates into improved performance. However, in treatments with a low bonus, initial orders are already too high, so focusing on the bonus simply causes orders to remain too high. In treatments without the bonus, orders start too high, but because there is no bonus on which to focus, they quickly go down.

When the review period is short, the probability of observing a 100% fill rate with a relatively low order

<sup>7</sup> For practical reasons, we could not offer the executives financial incentives. Instead, we asked them to help us with our research and give the decisions serious thought. Of the approximately 30 managers who initially agreed to participate, only eight completed both games. It is our opinion that those eight managers took the task seriously and the request to help us with our research was a sufficient incentive in this case.



is quite high. This high probability of 100% fill rate, combined with demand chasing, prevents the decision makers from learning to order enough in the high-bonus ( $B = 50$ ) condition. Because for high and medium bonuses participants start by placing orders that are below optimal, they need to learn to increase their orders, but in the high-bonus condition ( $B = 50$ ), this learning can only take place when the review period is sufficiently long.

The main practical implication from our study is that, in designing contracts based on fill-rate performance, managers should think carefully about the length of the review period. If the goal is to induce higher orders (e.g., use the SLA to address a double-marginalization problem), longer review periods may be more effective, even if the suppliers optimal order level is greater for a short horizon. Due to the high variability in feedback associated with short review horizons, a longer review period with a moderate bonus for achieving the target fill rate may well be more effective than a review period that is too short, even if it carries a higher bonus.

Order-up-to policies are not optimal in environments with SLAs. The optimal dynamic policies that adjust orders based on the fill rate achieved at any given point are complex. Analyzing those more complex policies, theoretically as well as in the laboratory, is a promising direction for future research. Another interesting issue that may be explored in the field is the extent to which managers faced with SLAs use static (and suboptimal) order-up-to policies or dynamic policies.

### Acknowledgments

The authors thank Axel Ockenfels and the Deutsche Forschungsgemeinschaft for financial support through the Leibniz-Program, the Smeal College of Business for support through the Smeal Summer Grants Program, and the Center for Supply Chain Research (CSCR) and the Smeal College of Business for their support. The authors also thank the executive subjects for their participation in the study and also the special issue editors, SE, and reviewers for their helpful feedback on this paper.

### References

- Banerjee, A., A. Paul. 2005. Average fill rate and horizon length. *Oper. Res. Lett.* **33**(5) 525–530.
- Bendoly, A., K. Donohue, K. Schultz. 2006. Behavior in operations management: Assessing recent findings and revisiting old assumptions. *J. Oper. Management* **24**(6) 747–752.
- Ben-Zion, U., Y. Cohen, R. Peled, T. Shavit. 2007. Decision-making and the newsvendor problem—An experimental study. Working paper. Ben-Gurion University of the Negev, Be'er-Sheva, Israel.
- Bolton, G., E. Katok. 2007. Learning-by-doing in the newsvendor problem: A laboratory investigation. *Manufacturing Service Oper. Management*. Forthcoming.
- Bostian, A. J., C. A. Holt, A. M. Smith. 2007. The newsvendor “pull-to-center effect”: Adaptive learning in a laboratory experiment. *Manufacturing Service Oper. Management*. Forthcoming.
- Cachon, G. 2003. Supply chain coordination with contracts. S. Graves, T. de Kok, eds. *Handbooks in Operations Research and Management Science: Supply Chain Management*. North Holland, Amsterdam, 229–340.
- Camerer, C. F. 2003. *Behavioral Game Theory Experiments in Strategic Interactions*. Princeton University Press, Princeton, NJ.
- Chen, J., D. Lin, D. Thomas. 2003. On the item fill-rate for a finite horizon. *Oper. Res. Lett.* **31**(2) 119–123.
- Goeree, J. K., C. A. Holt. 2005. An explanation of anomalous behavior in models of political participation. *Amer. Political Sci. Rev.* **99**(2) 201–213.
- Hogarth, R. 1987. *Judgement and Choice*, 2nd ed. John Wiley & Sons, New York.
- Kagel, J. H., A. E. Roth, eds. 1995. *Handbook of Experimental Economics*. Princeton University Press, Princeton, NJ.
- Kahneman, D., A. Tversky. 1979. Prospect theory: An analysis of decision under risk. *Econometrica* **47** 263–291.
- Kahneman, D., P. Slovic, A. Tversky. 1982. *Judgment Under Uncertainty: Heuristics and Biases*. Cambridge University Press, Cambridge, MA.
- Kay, E. 2005. Ways to measure supplier performance. Purchasing.com. <http://www.purchasing.com/article/CA508544.html?industryid=2161&nid=2419>.
- Keser, C., G. Paleologo. 2004. Experimental investigation of retailer-supplier contracts: The wholesale price contract. Working Paper 2004s-57, CIRANO, Montreal, Quebec, Canada.
- Lurie, N. H., J. M. Swaminathan. 2007. Is timely information always better? The effect of feedback frequency on performance and knowledge acquisition. Working paper, UNC Chapel Hill, Chapel Hill, NC.
- McKay, M. D., R. J. Beckman, W. J. Conover. 1979. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics* **21**(2) 239–245.
- Rapoport, A. 1966. A study of human control in a stochastic multistage decision task. *Behavioral Sci.* **11** 18–32.
- Rapoport, A. 1967. Variables affecting decisions in a multistage inventory task. *Behavioral Sci.* **12** 194–204.
- Schweitzer, M., G. Cachon. 2000. Decision bias in the newsvendor problem: Experimental evidence. *Management Sci.* **46**(3) 404–420.
- Selten, R., R. Stoeker. 1986. End behavior in sequences of finite prisoners' dilemma supergames: A learning theory approach. *J. Econom. Behav. Organ.* **7** 47–70.
- Spengler, J. J. 1950. Vertical integration and antitrust policy. *J. Political Econom.* **58**(4) 347–352.
- Taylor, S. E., S. T. Fiske. 1975. Point of view and perceptions of causality. *J. Personality Soc. Psych.* **32** 439–445.
- Thomas, D. 2005. Measuring item fill rate performance in a finite horizon. *Manufacturing Service Oper. Management* **7**(1) 74–80.
- Tversky, A., D. Kahneman. 1973. Availability: A heuristic for judging frequency and probability. *Cognitive Psych.* **5** 207–232.