

# Cross-Selling in a Call Center with a Heterogeneous Customer Population

Itay Gurvich\*

Mor Armony†

Constantinos Maglaras‡

August 29, 2006

## Abstract

Cross-selling is becoming an increasingly prevalent practice in call centers, due, in part, to its unique capability to allow firms to dynamically segment their callers and customize their product offerings accordingly. This paper considers a call center with cross-selling capability that serves a pool of customers that are differentiated in terms of their revenue potential and delay sensitivity. It studies the operational decisions of staffing, call routing, and cross-selling under various forms of customer segmentation. It derives near-optimal controls in each of the settings analyzed, and characterizes the impact of a more refined customer segmentation on the structure of these policies and the center's profitability.

## 1 Introduction

Many organizations consider their call centers as one of the most important channels of interaction with their customers, acting both as a service center and a point of sales – an opportunity for the firm to generate extra revenue by offering new or existing products to their customers. The significant revenue potential of this cross-selling strategy is underscored by the nature of the interaction that takes place in a call center and the wealth of information that is available through state-of-the-art Customer Relationship Management (CRM) systems; Together, they enable firms to segment their customer pools effectively and to tailor their product offerings to each such segment to increase the likelihood of purchase and the associated expected revenue. A familiar and successful example

---

\*Columbia Business School, 41 Uris Hall, 3022 Broadway, NY, NY 10027. (ig2126@columbia.edu)

†Stern School of Business, NYU, 44 West 4th Street, NY, NY 10012. (marmony@stern.nyu.edu)

‡Columbia Business School, 409 Uris Hall, 3022 Broadway, NY, NY 10027. (c.maglaras@gsb.columbia.edu)

of cross-selling practice is in the financial services industry, where customers that call for service, such as account balance inquiries, are often offered new financial products.<sup>1</sup>

Alongside its potential benefits, cross-selling may substantially increase the total workload that needs to be handled by the call center’s agents<sup>2</sup>. This may degrade the system’s quality of service, which is likely to have an adverse effect on the overall customer experience, as well as the effectiveness of cross-selling itself. Hence, it is important to carefully select which cross-selling opportunities to pursue and when to do so, and to account for the impact of these decisions in determining the staffing level of the call center. This paper considers a call center with cross-selling capabilities that serves a heterogenous pool of customers, and studies the operational decisions of staffing, call routing, and cross-selling under various forms of customer segmentation. It derives near-optimal controls in each of the settings analyzed, and characterizes the impact of more refined customer segmentation on the structure of these policies and the center’s profitability.

In more detail, we consider a call center with a single pool of fully flexible agents that first handle inbound call service requests, and subsequently decide if to attempt to cross-sell to some of these customers a certain product or service, whenever such an opportunity arises. Cross-selling attempts are handled by the same agent that has served the customer’s original request, upon completion of that task. Each cross-selling attempt is preceded by an instantaneous step that captures the customer’s decision of whether or not to agree to listen to the cross-selling offer. The processing times for the original service request and the cross-selling phase are exponentially distributed with potentially different parameters. Finally, the heterogenous pool of potential customers comprises of a discrete set of types or segments. Types differ in terms of their delay sensitivity and revenue potential. These are captured through the probability that a customer will agree to listen to a cross-selling offer as a function of the waiting time that he encountered, and through a demand relation that specifies the fraction of customers of a type that listened to the cross-selling offer and decided to buy the offered product as a function of the quoted price.

We study three variants of this model with an increasing degree of customer segmentation and, as a result, increasing flexibility in terms of the aforementioned operational and pricing decisions. The simplest model is one where customers are not segmented, or equivalently, where their types are not observable. In this case the manager is limited to make the cross-selling decisions based

---

<sup>1</sup>A recent study by McKinsey & Co. [11] suggests that bank call centers can generate through cross-selling revenues that are equivalent to 10% of the revenue generated through the retail branch channels.

<sup>2</sup>In a recent study, a Purdue University research group [2] estimates that call centers may attempt to cross-sell to as many as 60% of all its callers.

solely on the aggregate load in the system and to offer the same product to all customers. The second model is one where types are observed sometime during their service, and this information can therefore be used in deciding whether to cross-sell to a customer, and if so, which product to offer. In terms of product customization, we restrict attention to call centers that market the same product to all of its customers, e.g., a 6-month Certified Deposit (CD) account, but where the center has the capability to customize the price it quotes to customers of different types, e.g., by offering different interest rates for that CD. The third model is one where customer types are observable upon arrival, in which case the manager can also decide how to route customers of different types to the available agents. For each of these models, the call center manager’s problem is to select its staffing, routing, and cross-selling policies to maximize the center’s expected profit rate, given by its revenues minus the staffing cost minus a linear waiting time cost that is experienced by all customers and is incurred by the center.

The controlled two-stage service sequence of each customer and the dependence of the cross-selling phase on dynamic waiting time information makes an exact analysis of this model cumbersome and difficult, even if customers are treated as one segment. Our approach considers a deterministic relaxation of this problem, which is solved in closed-form. Its solution suggests different staffing and cross-selling policies for each of the model variants listed above. In each case, we show that our proposed policy is asymptotically optimal in systems with increasing call volume.

Our tractable deterministic analysis and the asymptotic performance guarantees of the proposed policies lead to several insights. The first one is that the marketing decisions of customer segmentation and product customization are effectively decoupled from the operational decisions of staffing, routing and cross-selling. Specifically, once the set of customer segments has been identified through an appropriate marketing and statistical analysis, and their respective characteristics have been identified using observed data<sup>3</sup>, the firm can precompute its product customization strategy ahead of time. The customized prices are then fed into the operational control problem that involves staffing, routing, and cross-selling decisions.

The degree of customer segmentation has many important consequences, which can also be easily seen from our deterministic relaxation. To start with, roughly speaking, the center will only

---

<sup>3</sup>The first step involves the identification of appropriate attributes along which to segment the customer pool. The accuracy of the estimation of the customer type characteristics will be greatly improved if the center can keep track of data on customers that refused to listen to the cross-selling offer, and on those that listened but did not buy. Finally, there is a tradeoff between the number of customer segments and the accuracy of this estimation procedure, which may result into coarse segmentation as opposed to segmenting down to the level of each customer.

cross-sell to customers that generate an expected revenue that exceeds the capacity cost involved in pursuing this attempt; the expected revenue is equal to the quoted price times the probability that this customer will buy the offered product. If the center can segment its customers, then it will only cross-sell to its profitable types; if no segmentation capability is in place, then it will either cross-sell to all customers or to none, depending again on the expected profitability of these cross-selling attempts. In each case, the center will staff so as to handle all regular service requests plus the additional nominal workload generated by its expected cross-selling activities. Since the cross-selling is controllable, it can provide enough flexibility in the use of the center's capacity, which eliminates the need to add "safety staffing" like is typically done according to the "square-root" rule in order to stabilize the system and guarantee moderate congestion. It is possible that even though it is profitable to cross-sell in a system that segments its customers, this is not the case without segmentation. Our analysis outlines such cases. Overall, customer segmentation increases the center's profitability in two ways: first, through a more efficient use of capacity achieved by reducing the volume of cross-selling attempts that are unlikely to be profitable, and second by customizing the product offering (price) for each customer type so as to maximize the resulting expected revenue. Finally, we note that the effect of observing the customer type upon arrival as opposed to after service has commenced is small. This is explained by the fact that even when the system does not differentiate between types in its routing decisions and handles all external calls through a common FCFS queue for all these types, the resulting waiting times are small; these are moderated through the dynamic cross-selling decisions of the call center and are reinforced by the customers' delay averseness.

The structure of the remainder of the paper is as follows. This section concludes with a brief literature survey. §2 describes the two models with observable types, emphasizing mostly the model where customer type is revealed once his service starts. These two models are analyzed in §3. §4 shows how the pricing problem can be treated separately from all other decisions, which is then used in §5 to analyze a model with no customer segmentation. §6 provides results from our numerical experiments and some concluding remarks, while the Appendix presents all of our proofs.

**Literature Review** The literature on the operational aspects of call centers is rich and has been increasing over the last decade. A survey of this literature and a tutorial on the subject can be found in [15]. Of particular relevance to our work is the literature on staffing of call centers. The most commonly used staffing rule in the literature is the so-called Square-Root Safety

Staffing rule, according to which the number of servers required to handle an offered-load of size  $R$  is  $R + \beta\sqrt{R}$ , for some constant  $\beta$ . The square-root safety staffing rule dates back to Erlang in his 1923 paper (that appeared in [12]). This rule was formalized by Halfin and Whitt [19] who showed that this square-root safety staffing rule guarantees very short delays in an appropriate asymptotic regime, and was shown to be nearly-optimal for a pure service center that handles a homogeneous customer population in Borst et. al. [8]. Finally, square-root safety staffing has been observed to be fairly robust with respect to changes in model assumptions to include features such as customer abandonment [16, 23], multiple customer classes [6, 5, 18], multiple server pools [3] and non-stationary arrival rates [13]. In contrast with the above stream of work, our paper shows that the issue of safety staffing is of lesser importance in call centers with significant cross-selling activity, since by adjusting the latter the manager can also control its congestion.

The cross-selling control problem has been studied by several authors, including Akşin and Harker [1], Örmeci and Akşin [24], and Byers and So [9]. These papers demonstrated some of the potential benefits of this practice, but also illustrated the high complexity of characterizing the optimal policy even when the staffing level is fixed and other simplifying assumptions are made.

To circumvent this issue, our paper focuses on large-scale systems and asymptotic optimality criteria that yield tractable problem formulations and lead to simple and practical control rules. This approach builds on earlier work in Armony and Gurvich [4], which to the best of our knowledge is the first paper to study the joint problem of staffing and control in a call center with a homogeneous customer population. Its main finding is that a threshold type cross-selling policy is nearly optimal. Our current paper extends the model of Armony and Gurvich [4] to consider a heterogeneous and delay sensitive customer population that is served by a call center that operates under varying forms of customer segmentation. This paper differs from [4] in its asymptotic framework, in the sense that we study call center systems of increasing size as measured by the nominal call volume, keeping all other model primitives such as the customer preferences and the cost structure unchanged; the latter were scaled in [4]. The deterministic relaxation that underlies our work is motivated by the work of Maglaras and Zeevi [21]. Finally, the economic model that we adopt and the notions of product differentiation and price discrimination that underlie our work are related to a vast literature in economics, marketing, and revenue management. We refer the reader to the book by Talluri and Van Ryzin [26] for an introduction of these issues.

## 2 Model formulation

We consider a call center with a single pool of  $N$  fully flexible agents that serves a heterogeneous customer population, comprising of  $K$  distinct segments, or types, or classes. We will study three model variants depending on the extent to which the customer types are observable by the system. These are graphically depicted in Figure 1. Model (a) assumes that types are unobservable, or that the call center does not segment its customers. In model (b), the type of a customer is observed when s/he is being served, and this information is subsequently used in the center’s cross-selling decisions. Finally, model (c) is one where the customer type is immediately observed upon arrival, e.g., by requiring customers to enter an account number, and can therefore be used in routing as well as in cross-selling decisions. We will focus on model (b), and treat model (c) as an extension and model (a) as a one-segment special case of this multi-segment model.

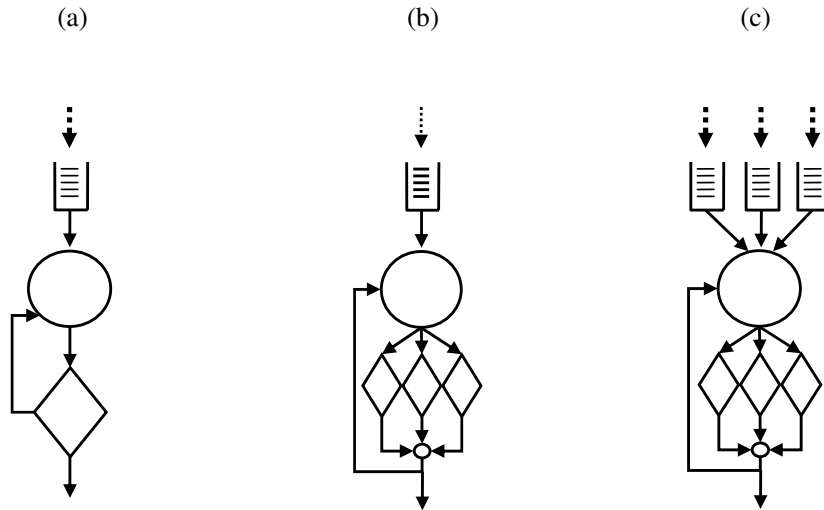


Figure 1: Three Cross-Selling Models

**Basic service:** Type  $i$  customers call the center according to a Poisson process,  $\{A_i(t), t \geq 0\}$ , with rate  $\lambda_i$ . Let  $A(t) = \sum_{i=1}^K A_i(t)$  and define  $\Lambda = \sum_{i=1}^K \lambda_i$  to be the total arrival rate into the system. The processing requirements are exponentially distributed with rate  $\mu_s$ , independent of the customer type. Under the assumption that types are unobservable before service begins (model (b)), all customers join a single queue and get processed in First-Come-First-Served (FCFS) manner.

**Cross-selling:** Once regular service is completed, a customer either leaves the system or enters

a cross-selling phase that is handled by the same agent. A cross-selling attempt is preceded by an instantaneous step whereat the customer is asked to listen to the actual offer, which itself requires an exponentially distributed amount of time with rate  $\mu_{cs}$ . All processing times (regular service and cross-selling) and inter-arrival times are assumed to be independent.

The probability that a type  $i$  customer will agree to listen to the cross-selling offer is given by  $q_i(w) = (q_i - a_i w)^+$ , where  $0 \leq q_i \leq 1$  and  $a_i > 0$  for  $i = 1, \dots, K$ , and  $w$  is the waiting time experienced by this customer before his service started;  $q_i < 1$  allows us to model cases where some customers may always decline to listen to the cross-selling offer.<sup>4</sup> If a class  $i$  customer agrees to listen to a cross-selling offer, he will be offered product  $i$  at price  $p_i$ . Class  $i$  customers have i.i.d valuations for this product, denoted by  $v_i$ , drawn from a distribution function  $F_i(\cdot)$ . Specifically, conditioned on agreeing to listen to a cross-selling offer, a class  $i$  customer will buy the product with probability  $\bar{F}_i(p_i)$  (where  $\bar{F}_i(\cdot) := 1 - F_i(\cdot)$ ), resulting in a conditional expected revenue from this customer of  $r_i := p_i \bar{F}_i(p_i)$ . We will also assume that  $p_i \bar{F}_i(p_i)$  are unimodal in the  $p_i$ 's for each  $i$ ; this is satisfied by many commonly used demand functions (see Talluri and van Ryzin [26]).

**Control decisions:** The call center manager selects the number of agents  $N$  in the system and has discretion with respect to the cross-selling decisions. We will consider policies,  $\pi$ , that decide whether to cross-sell to the  $j^{\text{th}}$  type  $i$  customer as a function of all the information available up to the decision point. In particular, the decision may depend on the customer's type, the waiting time encountered by this customer prior to his service, which we denote by  $w_{i,j}^\pi$ , the number of customers in the queue and the number of customers of each type  $i'$  that are currently in service, denoted by  $Q_{i'}^\pi(t)$  and  $Z_{i'}^\pi(t)$ , respectively. We let  $Q^\pi(t) = \sum_{i=1}^K Q_i^\pi(t)$  be the total queue length at time  $t$  under  $\pi$ . To guarantee the existence of steady state or at least the existence of long run averages for various quantities of interest we will restrict the set of admissible controls as follows.

**Definition 1 (*Admissible Controls*)** *Given a staffing level  $N$ , and parameters  $\lambda_1, \dots, \lambda_K, \mu_s, \mu_{cs}$ , we say that  $\pi$  is an admissible policy if it is non-preemptive, non-anticipative and  $\lim_{t \rightarrow \infty} E[Q^\pi(t)]/t \rightarrow 0$ . We denote the family of admissible policies by  $\mathcal{A}(\lambda_1, \dots, \lambda_K, \mu_s, \mu_{cs}, N)$ .*

Loosely speaking,  $\mathcal{A}(\lambda_1, \dots, \lambda_K, \mu_s, \mu_{cs}, N)$  is the set of stabilizing policies under the given parameters. Definition 1 takes into account the fact that the set of admissible policies depends on

---

<sup>4</sup>This model could arise if each type  $i$  customer is characterized through an i.i.d. "patience" parameter, such that he will accept to listen the cross-selling offer if his waiting time was less than his patience, and decline otherwise; the above model corresponds to a patience with a point mass of  $1 - q_i$  at 0, and is uniformly distributed in  $(0, q_i/a_i]$ .

the parameters of the model through the stability conditions of the system. To simplify notation, we will omit the the parameters  $\lambda_1, \dots, \lambda_K, \mu_s$  and  $\mu_{cs}$ , whenever these are exogenously fixed, and write  $\mathcal{A}(N)$  or simply  $\mathcal{A}$  whenever the staffing level is clear from the context. Note that the above definition implies that our system must be able to handle all of the nominal demand, at least when no cross-selling is exercised; that is, the staffing choice must satisfy the constraint  $N > R := \Lambda/\mu_s$ .

**Performance criterion:** We first define two system quantities that will play an important role in the call center's cost and revenue terms, respectively. Observe that a steady state need not exist for any  $\pi \in \mathcal{A}(N)$ . With that in mind, for any  $\pi \in \mathcal{A}(N)$  and  $i = 1, \dots, K$  we define

$$EW_i^\pi(\infty) := \limsup_{t \rightarrow \infty} \frac{\sum_{j=1}^{A_i(t)} w_{i,j}^\pi}{A_i(t)} \quad \text{and} \quad x_i(\pi) := \liminf_{t \rightarrow \infty} \frac{\sum_{j=1}^{A_i(t)} x_{i,j}^\pi}{A_i(t)},$$

where  $x_{i,j}^\pi$  is an indicator that is set to 1 whenever the  $j^{\text{th}}$  class  $i$  customer goes through a cross-selling phase, and  $x_{i,j}^\pi$  equals zero otherwise. When a steady state exists,  $EW_i^\pi(\infty)$  and  $x_i(\pi)$  coincide with the expected steady state waiting time experienced by type  $i$  customers, and the steady state fraction of class  $i$  customers that are asked *and* agree to listen to a cross-selling offer under  $\pi$ , respectively. Since customers are processed FCFS, it must be that  $EW_i^\pi(\infty) = EW_j^\pi(\infty)$  for all  $i, j$ , which will also be denoted by  $EW^\pi(\infty)$ .

The call center incurs linear staffing and waiting time costs per unit time, given by  $c \cdot N$  and  $\Lambda dEW^\pi(\infty)$ , respectively. The latter assumes that the waiting time cost is type independent. The waiting time cost can be thought of as a penalty that the system incurs in terms of lost goodwill from the customers. The call center manager's optimization problem is the following:

$$\sup_{N \in \mathbb{Z}_+, \pi \in \mathcal{A}(N)} \sum_{i=1}^K \lambda_i r_i x_i(\pi) - cN - \Lambda dEW^\pi(\infty). \quad (1)$$

Note that while it is not guaranteed that there exists a control that actually achieves the optimal profit rate, it is easy to establish the existence of an optimal  $N^*$ , since  $N$  is discrete, the profit rate is bounded above by  $\sum_i \lambda_i r_i - c \cdot R$ , and it decreases to  $-\infty$  as  $N$  grows large.

An alternate formulation to (1) would replace the waiting time cost by an upper bound constraint on the expected waiting time, typically 30 seconds, and consider the following problem:

$$\sup_{N \in \mathbb{Z}_+, \pi \in \mathcal{A}(N)} \left\{ \sum_{i=1}^K \lambda_i r_i x_i(\pi) - cN : EW^\pi(\infty) \leq \bar{W} \right\}. \quad (2)$$

Indeed, one can view (2) as a more natural starting point, and (1) as a “dualized” version of the problem that is perhaps simpler to address. We will refer to (1) and (2) as the **waiting cost** and **constrained** formulations, respectively. We will also make the following assumption:

**Assumption 1** *Types are labeled such that  $r_1 \geq \dots \geq r_K$  and  $r_1 > c/\mu_{cs}$ .*

The labeling assumption is innocuous. The condition  $r_1 > c/\mu_{cs}$  makes it profitable to cross-sell to at least type 1 customers. Without the latter, it would be economically optimal to restrict cross-selling to times when several agents are idle, which will tend to be infrequent if the call center staffing decision has itself been optimized.

### 3 Observable types: analysis based on a deterministic relaxation

A direct analysis of the problems formulated above is very difficult due to their multi-class nature and the dependence of the cross-selling success probability on state-dependent information. Our approach looks at relaxations of the above problems, where in addition to the staffing and cross-selling decisions, the manager can also select the waiting times experienced by its callers, which in reality are random variables that depend on the system dynamics. These relaxations are tractable, deterministic optimization problems that have insightful solutions and give rise to near-optimal heuristics. Focusing on model (b) (cf. Figure 1) first, §3.1 studies the waiting cost formulation of (1). These results are extended to the constrained formulation of (2) in §3.2, while §3.3 extends our work to model (c) where the customer types are observable upon arrival. All proofs are relegated to the appendix.

#### 3.1 The Waiting Cost Formulation

**Deterministic relaxation:** Starting with (1) we formulate the following linear program:

$$\begin{aligned}
 & \text{maximize} && \sum_{i=1}^K \lambda_i r_i x_i - c \cdot R(1+z) - d \sum_{i=1}^K \lambda_i w_i \\
 & \text{s.t.} && x_i \leq q_i(w_i), \quad \sum_{i=1}^K \lambda_i x_i \leq \mu_{cs} R z, \\
 & && z \geq 0; \quad x_i \geq 0, \quad w_i \geq 0, \quad \text{for all } i = 1, \dots, K,
 \end{aligned} \tag{3}$$

where  $x_i$  is interpreted as the fraction of class  $i$  customers that are being asked *and* agree to listen to a cross-selling offer;  $w_i$  is the “fictitious” waiting time experienced by class  $i$  customers in this

formulation; and  $z$  is the excess (normalized) staffing level beyond the nominal requirement of the offered load  $R$  ( $:= \Lambda/\mu_s$ ) as a fraction of  $R$ . The condition  $z \geq 0$  implies that the staffing level is sufficiently large to handle all basic service requests (i.e.,  $N \geq R$ ).

Recall the labeling convention in Assumption 1. Denoting the optimal solution to the knapsack problem in (3) with an overbar we have the following: set  $\bar{w}_i = 0$  for all  $i = 1, \dots, K$ ,

$$\bar{x}_i = \begin{cases} q_i & i \leq \bar{k} \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad \bar{z} = \sum_{i=1}^{\bar{k}} \frac{\lambda_i q_i}{R \mu_{cs}}, \quad (4)$$

where  $\bar{k} = \max\{i : r_i \geq c/\mu_{cs}\}$ . In fact, we will assume throughout that  $r_{\bar{k}} > c/\mu_{cs}$ , which is equivalent to assuming that the deterministic relaxation has a unique solution. In the presence of multiple solutions to the deterministic relaxation our approach might lead to multiple asymptotically optimal solutions. By Assumption 1,  $\bar{z}$  is guaranteed to be strictly positive. The resulting staffing level is  $R + \sum_{i=1}^{\bar{k}} \frac{\lambda_i q_i}{\mu_{cs}}$ . Note that the structure of the deterministic relaxation is such that as long as  $\lambda_i/\Lambda$  is known and is kept constant (which we will assume henceforth), the normalized quantities  $\bar{x}, \bar{z}$  do not change with  $\Lambda$ . Therefore, the relevant profit depends on the entire vector  $\lambda_1, \dots, \lambda_K$  through their sum  $\Lambda$  only. Specifically, the profit rate associated with the solution (4) is

$$\bar{\Pi}(\Lambda) = -cR + \sum_{i=1}^{\bar{k}} \lambda_i q_i (r_i - c/\mu_{cs}) = -cR + \sum_{i=1}^K \lambda_i q_i [(r_i - c/\mu_{cs}) \vee 0], \quad (5)$$

which is an upper bound for the optimal profit in (1). (Here and elsewhere  $x \vee y = \max\{x, y\}$ )

**A staffing and cross-selling proposal:** The nested structure of (4) is intuitive: we cross-sell to all types  $i$  for which their marginal revenue contribution,  $\lambda_i r_i q_i$ , exceeds the increase in staffing cost,  $c \frac{\lambda_i q_i}{\mu_{cs}}$ , resulting from the additional cross-selling workload; this reduces to the condition  $r_i > c/\mu_{cs}$ . The solution to the deterministic relaxation suggests the following pair of policies:

(S) *Staffing:* Staff with  $N = R(1 + \bar{z})$ .

(C) *Cross-selling:* Given the sequence of thresholds  $\eta_{\bar{k}} \leq \eta_{\bar{k}-1} \leq \dots \leq \eta_1$ :

- At any time  $t$ , define  $i(t) = \max\{i : \eta_i \geq Q(t)\}$ ;  $i(t) = 0$  if  $\eta_1 < Q(t)$ ;
- Cross-sell to any customer that completes service at time  $t$  and is of type  $i \leq i(t)$ .

The cross-selling policy (C) follows the solution of the deterministic relaxation when the queue length is modest, and then starts to reduce the amount of cross-selling activity as the system gets

increasingly congested. The asymptotic performance analysis that will follow does not use the precise values of the above thresholds, and in fact only makes use of the smallest threshold  $\eta_{\bar{k}}$ .

**Asymptotic optimality of (S)-(C):** Despite its simple structure, (S)-(C) performs very well in the stochastic system under consideration, and is, in fact, asymptotically optimal in large scale systems, i.e., where  $\Lambda$  is large. As a starting point we will establish that the system is always stable under (S)-(C) and that it admits a unique stationary distribution. We do that by showing the stronger result that the system will be stable under (C) as long as  $N > R$ , even if  $N < R(1 + \bar{z})$ .

**Proposition 1** *Fix  $\Lambda$  and assume (C) is used for some set of thresholds:  $\eta_{\bar{k}} \leq \eta_{\bar{k}-1} \leq \dots \leq \eta_1 \leq \infty$ . Then,  $N > R$  is a sufficient condition for stability. Moreover, for any  $N > R$ , the underlying Markov process admits a unique stationary distribution which is also its limiting distribution.*

This proposition illustrates this self-stabilizing nature of the cross-selling system. Note that the thresholds  $\eta_i$  are not needed for this result as they may all be set equal to  $\infty$ ; the stabilizing force stems from the delay sensitivity of the customers. Intuitively, when the system is heavily loaded, the queue and the resulting waiting time will grow large. In turn, fewer customers will agree to listen to cross-selling offers, thus reducing the load.

The remainder of this subsection will characterize the asymptotic performance of the original stochastic call center system under (S)-(C) in settings with large call volumes, as measured by  $\Lambda$ . One naturally expects that fixing a threshold policy, the best threshold values will be a function of the system size and in particular of  $\Lambda$ , the overall arrival rate. Let  $\eta_{\bar{k}}^\Lambda, \dots, \eta_1^\Lambda$ , be the threshold values corresponding to a system with arrival rate  $\Lambda$ . Then, we will show in our subsequent results that, indeed, there is a dependence of the threshold values on the system size and moreover that asymptotically optimal performance implies that these threshold values scale according to

$$\eta_i^\Lambda = \hat{\eta}_i \sqrt{\Lambda} \quad \text{for } i = 1, \dots, \bar{k} \quad (6)$$

and appropriate constants  $\hat{\eta}_{\bar{k}} \leq \dots \leq \hat{\eta}_1$ . Let  $N^*(\Lambda)$ ,  $x_i^*(\Lambda)$  and  $\Pi^*(\Lambda)$  denote the (unknown) optimal staffing level, realized long-run average cross-selling rates, and the corresponding profit rate for (1), respectively, when the aggregate demand is  $\Lambda$ . Also, let  $\hat{\Pi}(\Lambda)$  be the profit obtained when using (S)-(C) in the stochastic system. In the sequel we will make use of the following notation: for two positive sequences we say that  $x^\Lambda$  is  $o(y^\Lambda)$  if  $x^\Lambda/y^\Lambda \rightarrow 0$  as  $\Lambda \rightarrow \infty$ .

**Theorem 1** *Let  $\Lambda$  grow large, keeping  $\lambda_i/\Lambda$  constant for all  $i$ . Then, with thresholds satisfying (6), (S)-(C) is asymptotically optimal in the sense that*

$$\hat{\Pi}(\Lambda) = \Pi^*(\Lambda) - o(\Lambda). \quad (7)$$

Alternatively, one could write (7) in the form  $\hat{\Pi}(\Lambda)/\Pi^*(\Lambda) \rightarrow 1$  as  $\Lambda \rightarrow \infty$ . The proof of the above result follows by showing the stronger result that  $\hat{\Pi}(\Lambda)$  approaches  $\bar{\Pi}(\Lambda)$ , which itself is an upper bound for  $\Pi^*(\Lambda)$ . Since  $\Pi^*(\Lambda)$  is sandwiched between  $\hat{\Pi}(\Lambda)$  and  $\bar{\Pi}(\Lambda)$ , it must also be close to  $\bar{\Pi}(\Lambda)$ . This leads to a partial characterization of the unknown optimal policy in large scale systems.

**Theorem 2** *Let  $\Lambda$  grow large, keeping  $\lambda_i/\Lambda$  constant for all  $i$ . Then, i.  $\Pi^*(\Lambda) = \bar{\Pi}(\Lambda) - o(\Lambda)$ , ii.  $N^*(\Lambda) = R(1 + \bar{z}) \pm o(\Lambda)$ , and iii.  $x_i^*(\Lambda) = \bar{x}_i + o(1)$ .*

The above two Theorems together demonstrate how the solution of the deterministic relaxation captures the first order behavior of the optimal policy for (1), both in terms of its staffing and cross-selling decisions as well as its resulting profits. A closer look at the rate at which the waiting times converge to  $\bar{w} = 0$  will offer some insight on the selection of the threshold parameters  $\eta$  used in the cross-selling policy, as well as refine the above results. Specifically, the next lemma shows that if the thresholds  $\eta$  are of order  $\sqrt{\Lambda}$  (as in (6)), then the steady state waiting times that characterize the system are of order  $1/\sqrt{\Lambda}$ ; this is the nominal time it takes order  $\Lambda$  serves to clear a queue length of order  $\sqrt{\Lambda}$ . Thresholds of smaller magnitudes would result in even smaller waiting times.

**Lemma 1** *Let  $\Lambda$  grow large, keeping  $\lambda_i/\Lambda$  constant for all  $i$ . Denote by  $E[W^\Lambda]$  the steady-state expected waiting time under policy (S)-(C). Then, with thresholds satisfying (6),  $E[W^\Lambda] = O\left(1/\sqrt{\Lambda}\right)$ , or equivalently,  $\limsup_{\Lambda \rightarrow \infty} \sqrt{\Lambda}E[W^\Lambda] < \infty$ . In particular,  $E[W^\Lambda] \rightarrow 0$  as  $\Lambda \rightarrow \infty$ .*

Moreover, the next lemma shows that it is optimal to staff and cross-sell so that the waiting times are indeed of order  $1/\sqrt{\Lambda}$ , or smaller. We denote by  $E[W^{\Lambda,*}]$  the expected steady state waiting time under the optimal control  $(N^*(\Lambda), x^*(\Lambda))$ .

**Lemma 2** *Let  $\Lambda$  grow large, keeping  $\lambda_i/\Lambda$  constant for all  $i$ . If an optimal policy  $(N^*(\Lambda), x^*(\Lambda))$  exists, then  $\limsup_{\Lambda \rightarrow \infty} \sqrt{\Lambda}E[W^{\Lambda,*}] < \infty$ .*

Using the above results we can refine the conclusions of Theorem 2 to those given below:

**Proposition 2** *Let  $\Lambda$  grow large, keeping  $\lambda_i/\Lambda$  constant for all  $i$ . Then, i.  $\Pi^*(\Lambda) = \bar{\Pi}(\Lambda) - O(\sqrt{\Lambda})$ , ii.  $N^*(\Lambda) = R(1 + \bar{z}) \pm O(\sqrt{\Lambda})$ , and iii.  $x_i^*(\Lambda) = \bar{x}_i + O(1/\sqrt{\Lambda})$ .*

The values of the thresholds  $\eta_i$  can be selected via simulation. Our asymptotic analysis and our experience with numerical examples shows that only the smallest threshold,  $\eta_{\bar{k}}$ , has an important effect on the system performance; in non-stationary environments where  $\bar{k}$  may change as a function of the offered load, the larger thresholds may also play a role.

### 3.2 The Constrained Formulation

Lemmas 1 and 2 illustrate that the waiting times experienced in an optimally controlled call center will be of order  $1/\sqrt{\Lambda}$ . With that in mind, a waiting time constraint of the form  $E[W^\Lambda] \leq \bar{W}$  will become irrelevant as  $\Lambda$  grows large, since the actual waiting times will be much smaller than the desired target  $\bar{W}$ . A more appropriate formulation that would remain meaningful as  $\Lambda$  grows large would replace the upper bound constraint by a quantity that itself changes with  $\Lambda$  of the form  $\bar{W}^\Lambda = \hat{W}/\sqrt{\Lambda}$  for an appropriate choice of  $\hat{W}$ .<sup>5</sup> This would result in the following problem:

$$\sup_{N \in \mathbb{Z}_+, \pi \in \mathcal{A}(N)} \left\{ \sum_{i=1}^K \lambda_i r_i x_i(\pi) - cN : EW^\pi(\infty) \leq \bar{W}^\Lambda \right\}, \quad (8)$$

where  $\bar{W}^\Lambda = \hat{W}/\sqrt{\Lambda}$  for an appropriate choice of  $\hat{W}$ . Along the lines of (3) the following is a **deterministic relaxation** of (8):

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^K \lambda_i r_i x_i - c \cdot R(1 + z) \\ & \text{s.t.} && \sum_{i=1}^K \frac{\lambda_i}{\Lambda} w_i \leq \bar{W}^\Lambda, \\ & && x_i \leq q_i(w_i), \quad \sum_{i=1}^K \lambda_i x_i \leq \mu_{cs} R z, \\ & && x_i \geq 0, \quad w_i \geq 0, \quad \text{for all } i = 1, \dots, K. \end{aligned} \quad (9)$$

The linear program described above has the same optimal solution as (3), making our solution insensitive to the precise articulation of the effect of customer waiting times. The resulting staffing and cross-selling heuristics are again the ones described by (S)-(C) in the previous subsection. In

---

<sup>5</sup>For example, if the problem of original interest has  $\Lambda' = 100$  and  $\bar{W}' = 20$  seconds, then  $\hat{W}$  is selected so that  $\bar{W}' = \hat{W}/\sqrt{\Lambda'}$ , which in this case would give  $\hat{W} = 200$  seconds. One should then study an asymptotic version of (2) as  $\Lambda$  grows large and  $\bar{W}$  is scaled according to  $200/\sqrt{\Lambda}$ ; note that the original formulation is recovered for  $\Lambda = 100$ .

the case of the constrained formulation, one can also get a crude estimate for the threshold  $\eta_{\bar{k}}$  to be  $\eta_{\bar{k}} := \Lambda \bar{W}^\Lambda$ , which is consistent with (6). Intuitively, if the queue length is maintained below that threshold, then by a heuristic application of Little’s Law one would expect that the waiting times will be below  $\bar{W}^\Lambda$ . The next theorem establishes this result in an asymptotic sense as  $\Lambda$  grows large. With slight abuse of notation we use  $\hat{\Pi}(\Lambda)$  and  $\Pi^*(\Lambda)$  to denote the profit rate for the constrained formulation under (S)-(C) and the optimal policy, respectively.

**Theorem 3** *Let  $\Lambda$  grow large, keeping  $\lambda_i/\Lambda$  constant for all  $i$ . Then, with thresholds satisfying (6), i.  $\hat{\Pi}(\Lambda) = \Pi^*(\Lambda) + o(\Lambda)$  and ii.  $E[W^\Lambda] \leq \bar{W}^\Lambda + o(\bar{W}^\Lambda)$ .*

Theorem 3 shows that the waiting time constraint will be met asymptotically, but does not guarantee that it will be satisfied for any given value of  $\Lambda$ , including the one that characterizes the system of original interest. In order to ensure that the waiting time constraint will be met for a given system, e.g. with  $\Lambda' = 100$ , the manager should “fine tune” the queue length threshold  $\eta_{\bar{k}}$ . This can be done easily via simulation, or on-line through a feedback mechanism that adjusts the threshold based on the relation between the current average waiting time and its upper bound.

### 3.3 The value of customer type identification upon arrival

We complete the analysis of the model with observable types by comparing the model analyzed thus far (model (b) in Figure 1)) with the one where the type of each customer is observed at the time of his arrival to the system (model (c)). The latter could be achieved by requiring callers to identify themselves through a pin or an account number.

*Routing capability:* Once the call center observes the type of each arriving customer, it can maintain different (virtual) queues for customers of each type, and use that added flexibility in routing calls to available agents. This will eventually tradeoff the delay sensitivity and waiting time cost of each type against its potential revenue contribution. It is clear that this added element of control can only improve the call center’s profitability. The question is by how much. The main result of this section shows that the performance difference between FCFS routing (used when types were unobservable upon arrival) and any other routing policy that makes use of the type information, including the optimal one, is small and asymptotically negligible. The crude asymptotic analysis of this subsection uses a sandwich argument, similar to the one applied in Theorem 2, and does not need a detailed articulation of the set of admissible routing policies. We

refer the reader to Bassamboo *et al.*[7] for one possible definition of these controls.

Let  $\Pi^{**}(\Lambda)$  be the optimal achievable profit for the system where customer types are observable upon their arrival, and note that  $\Pi^{**}(\Lambda) \geq \Pi^*(\Lambda)$ . The key to our analysis is that the deterministic relaxations for models (b) and (c) are identical. The routing capability of model (c) can only serve to improve the vector of expected waiting times  $E[W_i]$ . Since the relaxation treats these as free optimization variables, denoted by  $w_i$ , and sets them equal to zero, its solution will coincide with that of (3). It follows that  $\bar{\Pi}(\Lambda) \geq \Pi^{**}(\Lambda) \geq \Pi^*(\Lambda)$ . From Theorem 2 we have that  $\Pi^*(\Lambda) = \bar{\Pi}(\Lambda) - o(\Lambda)$ , which leads to the following conclusion:

**Proposition 3** *Let  $\Lambda$  grow large, keeping  $\lambda_i/\Lambda$  constant for all  $i$ . Then *i.*  $\Pi^{**}(\Lambda) - \Pi^*(\Lambda) = o(\Lambda)$ . Specifically, *ii.*  $\Pi^{**}(\Lambda) - \Pi^*(\Lambda) = O(\sqrt{\Lambda})$ .*

So, while routing control capability may improve the quality of service enjoyed by some types and potentially simultaneously increase the revenue extracted from them, it will not lead to a significant overall profit gain. Moreover, the asymptotically optimal staffing and cross-selling recommendations that emerge from our analysis are insensitive (up to first order) to the use of this information.

The question that arises is whether segmentation at the cross-selling stage leads to significantly different results in comparison to no segmentation at all. To address this question we first study the issue of type-dependent product customization in §4, and then assess the value of customer segmentation in §5. For brevity, we will henceforth focus on a first order analysis in the spirit of Theorems 1 and 2, and drop the distinction between the waiting cost and constrained formulations.

## 4 The product customization problem

Customer segmentation in a call center setting allows firms to customize their products to better match the characteristics of each customer type and extract higher revenues. In our model, the firm can customize the price quoted to each customer type. In this section we show that the optimal prices can be computed separately from the operational decisions of staffing and cross-selling.

To this end, we will expand the notation used earlier on to let  $\Pi^*(\Lambda; p)$  and  $N^*(\Lambda; p)$  be the optimal profit rate and staffing level, respectively, for (1) for a given  $\Lambda$  and  $p$ . We then redefine  $\Pi^*(\Lambda) := \sup_{p \in \mathcal{P}} \Pi^*(\Lambda, p)$ , to be the optimal achievable profit rate when the call center is allowed to optimize over its price vector over some set  $\mathcal{P} = \mathcal{P}_1 \otimes \mathcal{P}_2 \otimes \dots \otimes \mathcal{P}_K$ , where for  $i = 1, \dots, K$ ,

$\mathcal{P}_i$  is assumed to be a compact interval in  $\mathbb{R}_+$ . Let  $p^*(\Lambda)$  be the optimal price vector, which is assumed to exist, and  $N^*(\Lambda)$  the corresponding staffing level. We also let  $\bar{\Pi}(\Lambda, p)$  be profit rate achieved in the deterministic relaxation of (3) for a given value of  $p$ ,  $\bar{\Pi}(\Lambda) := \max_{p \in \mathcal{P}} \bar{\Pi}(\Lambda, p)$ , be the profit rate when optimizing over the price, and let  $\bar{p}$  denote the corresponding optimizer, which may be different than  $p^*$ . While identifying  $p^*$  is hard, the deterministic price vector  $\bar{p}$  is easy to characterize by rewriting the objective function as

$$\bar{\Pi}(\Lambda, p) = -cR + \sum_{i=1}^K \lambda_i q_i [(r_i(p_i) - c/\mu_{cs}) \vee 0], \quad (10)$$

where  $r_i(p_i) = p_i \bar{F}_i(p_i)$ ; this expression reflects the fact that the center only cross-sells to and receives revenue from types for which  $r_i(p_i) \geq c/\mu_{cs}$ , and that it staffs accordingly. It follows that the corresponding optimal price is

$$\bar{p}_i = \operatorname{argmax}_{p_i \in \mathcal{P}_i} p_i \bar{F}_i(p_i), \quad (11)$$

and  $\bar{\Pi}(\Lambda) = -cR + \sum_{i=1}^K \lambda_i q_i [(r_i(\bar{p}_i) - c/\mu_{cs}) \vee 0] = \bar{\Pi}(\Lambda, \bar{p})$ . The corresponding staffing level is  $R(1 + \bar{z}(\bar{p}))$ , where

$$\bar{z}(\bar{p}) = \sum_{i=1}^{\bar{k}(\bar{p})} \frac{\lambda_i q_i}{R \mu_{cs}} \quad \text{and} \quad \bar{k}(\bar{p}) = \max\{i : r_i(\bar{p}_i) \geq c/\mu_{cs}\}; \quad (12)$$

the above expressions assume w.l.o.g that types are relabelled so that  $r_1(\bar{p}_1) \geq \dots \geq r_K(\bar{p}_K)$ . We also assume that  $r_i(\bar{p}_1) > c/\mu_{cs}$  and that  $r_{\bar{k}(\bar{p})}(\bar{p}_{\bar{k}(\bar{p})}) > c/\mu_{cs}$ , which guarantee, respectively, that Assumption 1 holds and that the solution of the deterministic relaxation given  $\bar{p}$  is unique. It is straightforward to show that  $\bar{p}$ ,  $\bar{z}(\bar{p})$  and  $\bar{k}(\bar{p})$  jointly characterize the optimal solution of the deterministic relaxation, and that this solution does not change if one were to scale  $\Lambda$  large, while keeping  $\lambda_i/\Lambda$  constant (this is the asymptotic setup adopted thus far). Note that although  $\bar{p}$  may be different than  $p^*(\Lambda)$ ,  $\bar{\Pi}(\Lambda, \bar{p})$  is still an upper bound for  $\Pi^*(\Lambda, p^*(\Lambda))$ . Using the results of the previous section we get the following:

**Proposition 4** *Define  $\bar{p}$ ,  $\bar{z}(\bar{p})$  through (11) and (12), respectively. Let  $\Lambda$  grow large, keeping  $\lambda_i/\Lambda$  constant for all  $i$ . Then: i.  $\Pi^*(\Lambda, p^*) = \bar{\Pi}(\Lambda, \bar{p}) - o(\Lambda)$ , ii.  $N^*(\Lambda, p^*) = R(1 + \bar{z}(\bar{p})) \pm o(\Lambda)$ , and iii.  $p^*(\Lambda) = \bar{p} + o(1)$ .*

An important consequence of the above result is that the pricing decisions can be done independently of the operational ones of staffing and cross-selling. This insight is valid in the system where types are observed upon arrival (model (c)), as well as in settings where products are customized along other non-price attributes which do not involve capacity and quality-of-service specifications.

## 5 The effect of customer segmentation

This section compares the profitability and behavior of the system studied in §3 and §4 against one that does not use a segmentation mechanism and instead treats its entire customer pool as one segment. The latter is offered a common product, i.e., at the same price, and cross-selling decisions are made without the customer type information; this is model (a) in Figure 1.

**A system with no customer segmentation:** The characteristics of this combined segment are a single delay sensitivity function  $q(\cdot)$  and a corresponding willingness-to-pay distribution  $F(\cdot)$  that are appropriate mixtures of the corresponding quantities for the various types. The delay sensitivity function,  $q(w)$ , is given by

$$q(w) := \sum_{i=1}^K \frac{\lambda_i}{\Lambda} q_i(w).$$

The combined willingness-to-pay distribution  $F$  is computed indirectly as follows. Let  $F(p|w)$  be equal to the probability that the willingness-to-pay of a customer that agreed to listen to the cross-selling offer after a waiting time of  $w$  time units is less than or equal to  $p$ . Then  $q(w)$  and  $\bar{F}(p|w)$  satisfy the following intuitive relation

$$q(w)\bar{F}(p|w) = \sum_{i=1}^K \frac{\lambda_i}{\Lambda} q_i(w)\bar{F}_i(p),$$

from which we can solve for  $F(p|w)$ .

The deterministic relaxation for the combined segment is now easy to solve by specializing the results of §3 to a single segment with characteristics  $q(w)$  and  $F(p|w)$ . Specifically, it is again optimal to set  $w = 0$ , which together with (10) gives the following objective

$$\bar{\Pi}^a(\Lambda, p) := -cR + \Lambda q(p\bar{F}(p|0) - c/\mu_{cs}) \vee 0, \quad (13)$$

where the superscript ‘a’ is meant to associate this expression to model (a), and

$$q := \sum_{i=1}^K \frac{\lambda_i}{\Lambda} q_i \quad \text{and} \quad \bar{F}(p|0) := \sum_{i=1}^K \frac{\lambda_i q_i}{\sum_{j=1}^K \lambda_j q_j} \bar{F}_i(p). \quad (14)$$

As shown in §4, one can study this deterministic formulation by separately optimizing over the price  $p$ , and then considering the resulting staffing and cross-selling problem at that price.

*The pricing decision:* The optimal price that the call center should use in this deterministic relaxation is given by the solution to the following problem:

$$\max_{p \in \bar{\mathcal{P}}} p \bar{F}(p|0), \quad (15)$$

which we denote by  $\bar{p}^a$ , and let  $r^a = \bar{p}^a \bar{F}(\bar{p}^a|0)$  and  $\bar{\mathcal{P}} = \mathcal{P}_1 \cap \mathcal{P}_2 \cap \dots \cap \mathcal{P}_K$ . Note that despite our assumptions regarding the unimodality of  $p_i \bar{F}_i(p_i)$ ,  $p \bar{F}(p|0)$  need not be unimodal itself. However, one can always find its optimizer through a single-parameter search.

*The staffing and cross-selling decisions:* Plugging  $\bar{p}^a$  into (13) and using the results of §3, the solution of the deterministic relaxation can be divided into two cases:

- i. If  $r^a \geq c/\mu_{cs}$ : the call center cross-sells to *all* customers and staffs with  $R_{\max} := R(1 + \bar{z}^a)$  servers, where  $\bar{z}^a = \Lambda q / (R \mu_{cs})$ .
- ii. If  $r^a < c/\mu_{cs}$ : the call center will *not* cross-sell to any customer and staff with  $R$  servers.

Using (13) and (14), the resulting profit rate in the deterministic relaxation is given by

$$\bar{\Pi}^a(\Lambda) := \begin{cases} -cR + \sum_{i=1}^K \lambda_i q_i (\bar{p}^a \bar{F}_i(\bar{p}^a) - c/\mu_{cs}), & \text{if } r^a > c/\mu_{cs} \\ -cR & \text{otherwise} \end{cases}, \quad (16)$$

which is again an upper bound for the optimal profit rate of the stochastic call center system.

As in §3, the natural implementation of the above policies in case i. would be to cross-sell as long as the queue is below an appropriate threshold that serves to limit excessive delays. In case ii., the system may still elect to cross-sell, but only if either the queue is very small or there are a sufficient number of agents that are idle. Moreover, in that case the staffing level should be inflated to  $R + x\sqrt{R}$  for an appropriate constant  $x$ , in order to provide moderate delays. The asymptotic analysis of §3 does apply to the single-segment model when the solution of the

deterministic relaxation falls into case i., but it does not cover case ii., where the system exercises negligible cross-selling. That case was studied in detail in Armony and Gurvich [4] and will not be further reviewed here.

**The effect of customer segmentation:** The key differences between the two systems, with and without segmentation, are best illustrated through their respective deterministic relaxations, which are simple and accurate, in the sense that they capture the structure of the underlying optimal policies and their resulting performance asymptotically.

1. *Cross-selling (all-or-none vs. selected types):* For both models, the call center will do significant cross-selling only if the expected revenue from doing so exceeds the capacity cost involved in that activity. With no segmentation capability in place, the system will either choose to cross-sell to all of its callers if  $r^a \geq c/\mu_{cs}$ , or to none. In the first case, this may involve cross-selling to customer segments to which it is strictly unprofitable to do so, while in the second it does involve foregoing profitable cross-selling opportunities that cannot be singled out from the larger pool of callers (the latter follows from Assumption 1). Using customer segmentation, the system will only cross-sell to types  $i = 1, \dots, \bar{k}$  for which  $r_i(\bar{p}_i) \geq c/\mu_{cs}$ , i.e., for which cross-selling is profitable. Finally, we note that although Assumption 1 guarantees that the call center will always choose to cross-sell to some subset of the customer types, if these can be segmented out, it does not guarantee that it is profitable to do so in a system with no segmentation capability.

2. *Staffing:* The model with no segmentation will either staff with  $R_{\max} = R(1 + \bar{z}^a)$  or  $R + x\sqrt{R}$  servers, depending on whether it will cross-sell or not. In contrast, the model with segmentation will staff with  $R(1 + \bar{z})$  servers;  $\bar{z} < z^a$ , unless it is profitable to cross-sell to all customer types.

3. *Uniform vs. customized pricing & profitability:* Most structural differences between the two systems originate from the pricing policies adopted by the call center in each case, and the corresponding expected revenue that they will generate per customer that agrees to enter the cross-selling phase. As explained earlier, the system that segments its customers will customize its prices  $\bar{p}_i$  for each type according to (11), while the system with no segmentation will use a uniform price,  $\bar{p}^a$  defined through (15). An immediate consequence of the above is that

$$r^a = \sum_{i=1}^K \frac{\lambda_i q_i}{\sum_{j=1}^K \lambda_j q_j} \bar{p}^a \bar{F}_i(\bar{p}^a) \leq \sum_{i=1}^K \frac{\lambda_i q_i}{\sum_{j=1}^K \lambda_j q_j} \bar{p}_i \bar{F}_i(\bar{p}_i).$$

Pre-multiplying by  $\sum_{j=1}^K \lambda_j q_j$  and subtracting out the corresponding capacity cost we get that

$$\left(\sum_{j=1}^K \lambda_j q_j\right) (r^a - c/\mu_{cs}) \leq \sum_{i=1}^K \lambda_i q_i (r_i(\bar{p}_i) - c/\mu_{cs}) \leq \sum_{i=1}^K \lambda_i q_i (r_i(\bar{p}_i) - c/\mu_{cs})^+.$$

The right-hand-side (RHS) of the above expression is equal to the profit contribution due to cross-selling in the system with segmentation, which is clearly non-negative. This allows us to strengthen this inequality to the following

$$\left(\sum_{j=1}^K \lambda_j q_j\right) (r^a - c/\mu_{cs})^+ \leq \sum_{i=1}^K \lambda_i q_i (r_i(\bar{p}_i) - c/\mu_{cs})^+, \quad (17)$$

where, in turn, the left-hand-side (LHS) of (17) is the profit contribution due to cross-selling in the system with no segmentation. The above inequality is strict as long as there exists a type  $i$  for which  $\bar{p}^a \bar{F}_i(\bar{p}^a) < \bar{p}_i \bar{F}_i(\bar{p}_i)$ , which by the definition of  $\bar{p}_i$  and the unimodality of  $p \bar{F}_i(p)$ , reduces to

$$\exists i \in \{1, \dots, K\} \text{ for which } \bar{p}_i \neq \bar{p}^a, \quad (18)$$

or equivalently to

$$\exists i, j \in \{1, \dots, K\}, \text{ such that } \bar{p}_i \neq \bar{p}_j. \quad (19)$$

Unless customer types have trivial differences with respect to their willingness-to-pay, conditions (18) or (19) are likely to be satisfied, in which case the ability to segment the customer pool would lead to significant profit gains. For example, if the willingness-to-pay distributions for the various types were exponential with parameters  $b_i$ , then the above conditions would require that at least two of these types had different parameters  $b_i \neq b_j$ . If the distributions were logistic with scale parameters  $b_i$  (these are commonly used in the literature in modelling different customer segments), then again (18) would require that the parameters of at least two segments are different. A simple extension of our previous results yields the following characterization of the potential value of customer segmentation in the underlying stochastic call center systems.

**Proposition 5** *Under Assumption 1, if (18) (or equivalently (19)) holds, then for all  $\Lambda$ ,  $\bar{\Pi}(\Lambda) - \bar{\Pi}^a(\Lambda) = \delta\Lambda$ , where  $\delta$  is the difference of the RHS and LHS of (17) normalized by  $\Lambda$ . Moreover, if we let  $\Lambda$  grow large, keeping  $\lambda_i/\Lambda$  constant for all  $i$ , then*

$$\Pi^*(\Lambda) - \Pi^{*,a}(\Lambda) = \delta\Lambda + o(\Lambda),$$

where  $\Pi^*(\Lambda)$ ,  $\Pi^{*,a}(\Lambda)$  are the optimal expected profit rates for the underlying stochastic systems with and without segmentation, respectively.

The above proposition together with the results of Theorems 1 and 3 suggest that the staffing and cross-selling policies proposed in this paper would realize most of the profit differential that can be attributed to customer segmentation. Operationally, the latter also leads to more efficient capacity utilization since call centers that do not segment their callers but try to cross-sell to them, end up pursuing too many customer prospects that are unlikely to lead to a sale. Our stylized, yet insightful, analysis can be used to assess the magnitude of this potential benefit, which is useful in deciding the value proposition of an investment in technology and agent training that would be needed to support a sophisticated customer segmentation and cross-selling strategy.

## 6 Numerical results and concluding remarks

This section first describes a set of numerical experiments that illustrate some of the key findings of this paper, and then closes with some concluding remarks.

### 6.1 Numerical experiments

Our results are organized in three categories. The first offers a representative numerical illustration of the accuracy of our asymptotic analysis. The second examines the quality of the proposed policies, and in particular shows the sensitivity of the system performance to changes in staffing and threshold levels that are used in the cross-selling decisions. The last one gives some examples of the potential value of using customer segmentation in such a call center.

**The accuracy of large-scale asymptotics:** We illustrate the accuracy of the proposed (S)-(C) heuristic by experimenting on a system with 4 customer classes. The service rates are  $\mu_s = 1$  and  $\mu_{cs} = 2$ ; one may regard all subsequent parameters as normalized with respect to  $\mu_s$ . The arrival rates are  $\lambda_1 = \lambda_2 = \frac{1}{3}\Lambda$  and  $\lambda_3 = \lambda_4 = \frac{1}{6}\Lambda$ , while the aggregate arrival rate,  $\Lambda$ , will be varied over a range of values in our experiment. The product prices are exogenously given and result in expected revenues per type  $i$  customer that goes through cross-selling given by  $r_1 = 7$ ,  $r_2 = 5$ , and  $r_3 = r_4 = 0.4$ . For simplicity we assume that the customers' delay sensitivity parameters are common across types and given by  $q_i = 1$  and  $a_i = 0.1$ . The staffing cost is normalized to  $c = 1$

and for concreteness we consider the constrained formulation with an upper bound for the waiting time equal to  $1/6$ ; if the natural time units are minutes, then this upper bound is 10 seconds. Under this choice of parameters we have that  $\bar{z} = \frac{1}{3} > 0$  and  $\bar{k} = 2$ , i.e., the center will cross-sell to types 1 and 2 only. These values of  $\bar{z}$  and  $\bar{k}$  and the above set of revenue and cost parameters give  $\bar{\Pi}(\Lambda) = (2.67)\Lambda$  as an upper bound on the system's profit rate.

We have simulated the system behavior under three variants of the policy (S)-(C) for  $\Lambda$  ranging from 40 to 200. The first variant is a direct translation of the solution of the deterministic relaxation, with a threshold  $\eta_2 = \lceil \frac{1}{6}\Lambda \rceil$ ; recall that type 2 is the least profitable type that the system cross-sells to. (For simplicity we set  $\eta_1 = \infty$ , i.e., the system would always cross-sell to type 1 customers.) The other two policy variants had  $\eta_2$  and the staffing level  $N$  further optimized via exhaustive simulation. The simulation code was written in c++. Each sample path contained 800,000 customer arrivals from which we formed time averages of the queue length and of the fraction of customers of each type that were cross-sold to. The length of each simulated path ensured that our estimates were close to the actual steady state behavior.

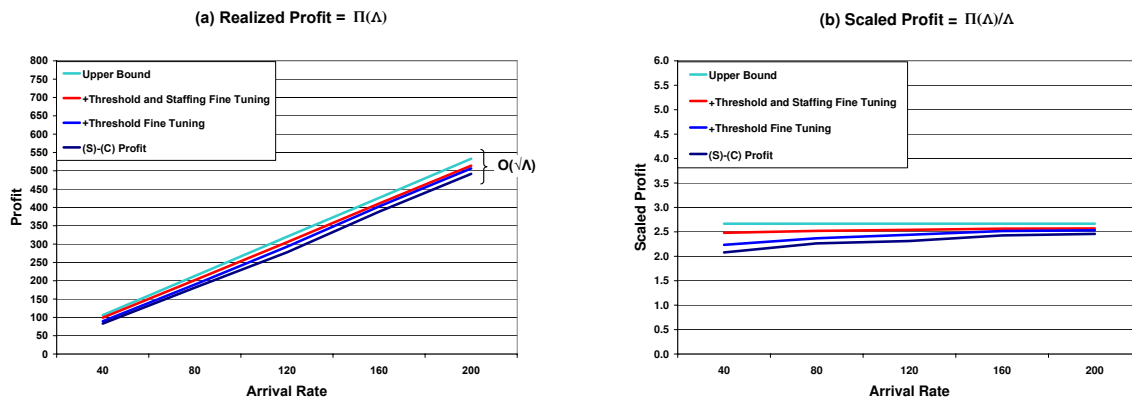


Figure 2: Performance of (S)-(C)

First, we note from Figure 2(a) that the absolute deviation between the profits achieved through the three candidate policies as well as their difference against the deterministic upper bound increases with the scale of the system, as measured by the aggregate call volume  $\Lambda$ . However, Figure 2(b) illustrates that if normalized by  $\Lambda$ , which is the multiplicative factor by which the above quantities are growing, then the respective difference decays to zero. In fact, this decay is of order  $1/\sqrt{\Lambda}$ . The above findings are representative of many examples that we tested. Second, we observe that as the size of the system increases, most of the profit gains from fine-tuning the cross-selling

threshold parameter and staffing level can be attributed to the former. This is practically appealing as it makes the model more robust to forecasting errors, because adjustments can be made on-line. The next set of results that we present study this issue in more detail, and also review the waiting time constraint qualification.

**Performance sensitivity with respect to the cross-selling threshold and the staffing level:** Figure 3 offers a more detailed look at the effect of these two parameters to the center’s profitability and the steady-state expected waiting time experienced by its callers for the system examined above for  $\Lambda = 120$ . The parameters extracted from the deterministic relaxation are  $\bar{z} = 1/3$  and  $\bar{k} = 2$ , which would translate to a nominal staffing of  $N = 160$  servers, and a nominal threshold of  $\eta_2 = \Lambda \bar{W}^\Lambda = 20$ ; i.e., the center would stop cross-selling to type 2 customers when there are more than 20 customers in queue. Specifically, Figure 3(a) shows the distance between the realized profit and its upper bound for various values of  $\eta_2$  and  $N$ . Figure 3(b) depicts the expected waiting time for each of these parameter combinations; the respective constraint requires that this falls below 0.167.

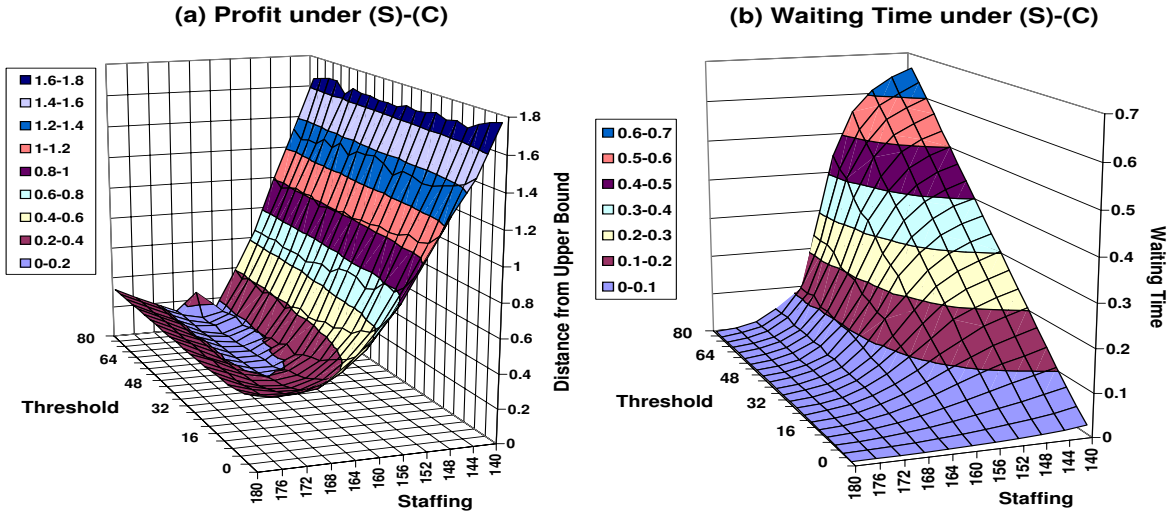


Figure 3: Performance as function of staffing and threshold levels

It is worth noting that the center’s profitability is fairly insensitive to the staffing level around its nominal value of 160 servers, since the effect of the latter can be compensated by appropriately adjusting the cross-selling threshold. As expected the waiting time is decreasing in the staffing level, and increasing in the value of the cross-selling threshold; i.e., more servers reduce the overall load,

while higher thresholds imply that the system is willing to tolerate longer waiting times. In fact, as expected from an informal application of Little’s law, the expected waiting time increases almost linearly as a function of the threshold. The effect of the threshold on the profit is less significant, which is consistent with our asymptotic results that showed that (S)-(C) with practically any threshold level performs very close to the upper bound in large systems. Taken together, the above comments suggest that call centers of reasonably large size can use the nominal staffing level extracted through the deterministic analysis, and subsequently select the cross-selling threshold to achieve constraint qualification and improve profits.

**The value of market segmentation:** We conclude this section through a set of numerical experiments that strive to illustrate the potential value of market segmentation. The analysis here is crude in the sense that it is limited to the deterministic relaxation. The asymptotic performance guarantees and the numerical results presented above suggest that the profit gap between the respective deterministic relaxations will persist in the stochastic systems as well. To facilitate the presentation of our results we will mostly focus on a two type system, for which  $\mu_s = 1$ ,  $\mu_{cs} = 2$ ,  $c = 1$ ,  $\Lambda = 100$ , and  $\lambda_1 = \lambda_2 = 0.5\Lambda$ . The waiting cost  $d$  or the waiting time upper bound  $\bar{W}$  do not play any role in the deterministic analysis, and hence there is no need to specify them.

It remains to specify the customer choice behavior. As in the previous examples, we assume that the delay preferences of both types are the same with  $q_i = 1$  and  $a_i = 0.1$  for  $i = 1, 2$ . Type  $i$  customers are assumed to have an exponentially distributed willingness-to-pay with parameter  $b_i$  for which  $\bar{F}_i(p_i) = e^{-b_i p_i}$ ,  $i = 1, 2$ . We assume that prices can obtain values on the bounded interval  $[0, 20]$  in each case. For the system that segments its customers, the optimal prices are given by  $\bar{p}_i = \frac{1}{b_i} \wedge 20$  (where  $x \wedge y = \min\{x, y\}$ ), for which  $r_i(\bar{p}_i) = \frac{1}{b_i} \wedge 20e^{-(20b_i \wedge 1)}$ . Note that the optimal price  $1/b_i$  in the absence of the price bound of \$20 is equal to the average of the distribution  $F_i$ , and that  $r_i(\bar{p}_i)$  is linear in  $1/b_i$  as long as  $b_i \geq 0.05$ . The solution to the deterministic relaxation will cross-sell to type  $i$  provided that  $r_i(\bar{p}_i) \geq c/\mu_{cs}$ , which in this model translates to  $b_i \leq 0.74 (= 2/e)$  and that  $1/b_i \geq 1.36$ . The optimal price for the system that cannot segment the two customer types does not admit a closed form solution, and is computed numerically using (14) and (15).

To test specific numerical system instances we have generated 250 independent realizations of the pair  $(b_1, b_2)$  by drawing each of the  $b_i$ ’s independently from a uniform distribution on  $[0, 2]$ . For each realization of  $(b_1, b_2)$  we solved the deterministic relaxations with and without segmentation. This involved computing the optimal prices, deciding to which types to cross-sell to, if any, calculating

the corresponding the staffing level, and finally the profit rate. Figure 4 displays the relative increase in profits,  $(\bar{\Pi}(\Lambda) - \bar{\Pi}^a(\Lambda))/\bar{\Pi}^a(\Lambda)$ , versus the maximum of the average willingness-to-pay among the two types, given by  $\max\left(\frac{1}{b_1}, \frac{1}{b_2}\right)$ . The average profit increase through segmentation in this two class experiment was around 24%. We have repeated this experiment several times and in all of the experiments the average profit was above 20%.

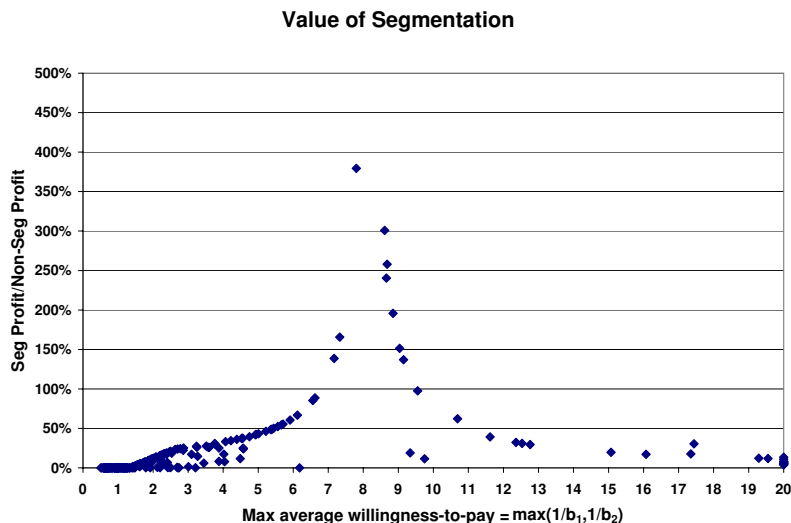


Figure 4: Profit comparison of systems with and without customer segmentation

Figure 4 is rather intuitive. There will be no profit gap between the two systems if  $b_1 = b_2$  or if the  $b_i$ 's are different but are such that no system decides to cross-sell to any customer. In settings where at least one type has a very large average willingness-to-pay, both systems will be very profitable in their cross-selling activities, and the relative difference in profit will be small (the RHS of the figure). In settings where both parameters  $1/b_i$  are small, then again the profit differential will be small because cross-selling is barely compensating for the cost of capacity. The difference between the two systems is more pronounced when  $\frac{1}{b_1}$  and  $\frac{1}{b_2}$  are of moderate size, in which case the relative added value from a) price customization and b) selective cross-selling (i.e., the capability to cross-sell to only one of the two types) is significant. For example, 20% of the 250 instances that we generated are such that the system with segmentation will choose to only cross-sell to one type, whereas the system with no segmentation capability will not cross-sell at all.

Finally, as the number of customer types and the degree of potential segmentation increases, the overall profit contribution due to segmentation becomes more substantial. In a set of experiments

that we ran with four customer types the average relative profit increase was 40% (up from 24% for the system with 2 types). Also, as the number of types was increased, we observed more instances where the cross-selling recommendations of the two systems would differ significantly.

## 6.2 Concluding Remarks

To summarize, this paper proposes a tractable deterministic relaxation for studying the various control problems that arise in call center systems with cross-selling capability, paying particular attention to the effect of customer segmentation on the structure of the staffing, cross-selling, and routing policies that the system may choose to adopt. The policies that are generated through this analysis are simple to implement, intuitive, and achieve near-optimal performance.

Our analysis can be extended in several directions to better model the operational complexity of modern call center systems, as well as that of customer behavior. In the former, this may include systems that have multiple pools of agents with different processing capabilities, as well as more complicated service requirements, that may need a sequence of steps to be handled by the same or different agents. With respect to the latter, one could model non-linear delay sensitivity functions, i.e., of the form  $q_i(w) = q_i f_i(w)$ , for continuous, strictly decreasing functions  $f_i(\cdot)$ , or allow the customers decision of whether to listen to the cross-selling offer to include information from the initial phase of service experienced by the customer, such as his service time, whether his initial request was successfully resolved, etc. Another extension would be to allow for customers to abandon the queue if their waiting time is excessive. All of the above generalizations increase the complexity of the underlying system substantially, but can be addressed using our approximate analysis with little additional effort. Another interesting extension would demonstrate the robustness properties of systems with cross-selling capability against non-stationary arrival patterns and parameter estimation and forecasting errors, as well as their advantage in producing more practical staffing schedules when faced with these operating challenges.

## References

- [1] O. Z. Akşin and P. T. Harker. To sell or not to sell: Determining the trade-offs between service and sales in retail banking phone centers. *Journal of Service Research*, 2(1):19–33, 1999.
- [2] J. Anton. Best practices in cross-selling and up-selling. *BenchmarkPortal.com*, 2005.

- [3] M. Armony. Dynamic routing in large-scale service systems with heterogeneous servers. *Queueing Systems*, 51(3-4):287–329, 2005.
- [4] M. Armony and I. Gurvich. When promotions meet operations: Cross-selling and its effect on call-center performance. *Working Paper*.
- [5] M. Armony and C. Maglaras. Contact centers with a call-back option and real-time delay information. *Oper. Res.*, 52(4):527–545, 2004.
- [6] M. Armony and C. Maglaras. On customer contact centers with a call-back option: customer decisions, routing rules and system design. *Oper. Res.*, 52(2):271–292, 2004.
- [7] A. Bassamboo, J. Harrison, and A. Zeevi. Design and control of a large call center: Asymptotic analysis of an lp-based method. *Operations Research*, 54:419–435, 2006.
- [8] S. Borst, A. Mandelbaum, and M. Reiman. Dimensioning large call centers. *Operations Research*, 52(1):17–34, 2004.
- [9] R. E. Byers and R. So. The value of information-based cross-sales policies in telephone service centers. *Working Paper*.
- [10] J. G. Dai. On positive Harris recurrence of multiclass queueing networks: A unified approach via fluid limit models. *Ann. Appl. Prob.*, 5:49–77, 1995.
- [11] A. Eichfeld, T. Morse, and K. Scott. Using call centers to boost revenue. *McKinsey Quarterly*, May 2006.
- [12] A. Erlang. On the rational determination of the number of circuits. In E. Brockmeyer, H. Halstrom, and A. Jensen, editors, *The life and works of A.K. Erlang*. Copenhagen: the Copenhagen Telephone Company, 1948.
- [13] Z. Feldman, A. Mandelbaum, W. Massey, and W. Whitt. Staffing of time-varying queues to achieve time-stable performance. 2005. Working Paper.
- [14] D. Gamarnik and A. Zeevi. Validity of heavy traffic steady-state approximations in generalized jackson networks. *Ann. Appl. Prob.*, 16:56–90, 2006.
- [15] N. Gans, G. Koole, and A. Mandelbaum. Telephone call centers: Tutorial, review, and research prospects. *Manufacturing & Service Operations Management*, 5(2):79–141, 2003.

- [16] O. Garnett, A. Mandelbaum, and M. Reiman. Designing a call center with impatient customers. *Manufacturing & Service Operations Management*, 4(3):208–227, 2002.
- [17] I. Gurvich, M. Armony, and C. Maglaras. Cross-selling in call centers with heterogeneous customer populations: Technical appendix.
- [18] I. Gurvich, M. Armony, and A. Mandelbaum. Service level differentiation in call centers with fully flexible servers.
- [19] S. Halfin and W. Whitt. Heavy-traffic limits for queues with many exponential servers. *Oper. Res.*, 29(3):567–588, 1981.
- [20] I. Karatzas and S. Shreve. *Brownian Motion and Stochastic Calculus*. 2nd ed., Springer-Verlag, New York, 1991.
- [21] C. Maglaras and A. Zeevi. Pricing and design of differentiated services: Approximate analysis and structural insights. *Oper. Res.*, 53(2):242–262, 2005.
- [22] A. Mandelbaum, W. Massey, and M. Reiman. Strong approximations for Markovian service networks. *Queueing Systems*, 30:149–201, 1998.
- [23] A. Mandelbaum and S. Zeltyn. Staffing many-server queues with impatient customers: constraint satisfaction in call centers. *Submitted*, 2006.
- [24] E. L. Örmeci and O. Z. Akşin. Revenue management through dynamic cross-selling in call centers. *Working Paper*.
- [25] A. Puhalskii. On the invariance principle for the first passage time. *Math. Oper. Res.*, 19:946–954, 1994.
- [26] K. Talluri and G. van Ryzin. *The theory and practice of revenue management*. Kluwer Academic Publishers, 2004.

## A Proofs

This section is dedicated to performance analysis of the (S)-(C) staffing and control rule. The analysis, which consists of several components, will eventually lead to the asymptotic optimality results of Theorems 1, 2 and 3 as well as Proposition 2. Some of the results proved below are based

partially on auxiliary results whose proof is involved and very lengthy. Hence, while we do give here the proofs for all the main results given in the paper, the proofs for some of auxiliary results are relegated to an online appendix [17].

The following notational conventions will be used throughout this appendix.  $Z^\Lambda(t)$  is the number of customers receiving service (not cross-selling), at time  $t$ . Recall that  $Q^\Lambda(t)$  is the overall queue length at time  $t$ . Also, let  $Z_{i,2}^\Lambda(t)$  be the number of class  $i$  customers in the cross-selling phase at time  $t$  and  $Z_2^\Lambda(t) = \sum_{i=1}^K Z_{i,2}^\Lambda(t)$  be the overall number of customers in the cross-selling phase at time  $t$ . We assume that all processes and random variables are defined on a common probability space  $\{\Omega, \mathcal{F}, P\}$  on which we make later additional probabilistic assumptions. For any finite dimensional random variable  $X$ ,  $|X|$  is the  $L^1$  norm. For any finite dimensional process  $\{Y(t), t \geq 0\}$ ,  $E_y[|Y(t)|]$  stands for the expected value of the normed process at time  $t$  given that  $Y(0) = y$ .

### Proof of Proposition 1:

First note that the delay sensitivity of the customers dictates that we need to keep track of the individual customers' waiting times to generate the sample path of the system. Hence, the state descriptor  $S^\Lambda(t) = (Z^\Lambda(t); Z_{i,2}^\Lambda(t), i = 1, \dots, K; Q^\Lambda(t))$  is not large enough for a Markovian characterization under (S)-(C). Instead, we use a larger state descriptor that contains additional workload information. Specifically, let  $v_j^\Lambda(t)$  be the residual handling time of the  $j^{\text{th}}$  arriving customer (service and cross-selling) at time  $t$ . Then, we construct the vector  $v^\Lambda(t) = \{v_j^\Lambda(t)\}_{j \geq 1: v_j^\Lambda(t) > 0}$ , and consider the process  $\Xi^\Lambda(t) = \{v^\Lambda(t), Q^\Lambda(t)\}$  so that for all  $t \geq 0$ ,  $\Xi^\Lambda(t) \in \Xi := \mathbb{R}_+^\infty \times \mathbb{Z}_+$ .

The sample paths of the process  $\Xi^\Lambda(t)$  are generated as follows: we generate an infinite sequence of I.I.D uniform  $[0, 1]$  random variables and let customer  $j$  agree to cross-selling if  $U_j \leq q(W(\tau_j))$  where  $\tau_j$  is the arrival time of customer  $j$ . We generate an infinite sequence  $\{s_j^1, s_j^2\}_{j=1}^\infty$  of I.I.D random variables, where for each  $j$ ,  $s_j^1$  and  $s_j^2$  are independent exponentially distributed random variables with respective rates  $\mu_s$  and  $\mu_{cs}$ . We let  $s_j(w^j)$  be the handling time of customer  $j$  given that he had to wait  $w^j$  units of time. Then, if the agent decided to cross-sell to customer  $j$ ,  $s_j(w^j)$  will equal  $s_j^1 + s_j^2$  with probability  $q(w^j)$  and it will equal  $s_j^1$  otherwise. Note, that under (S)-(C) and with the given state descriptor one can calculate the waiting time (and in turn the actual service time) of customer  $j$  immediately upon the customer's arrival to the system. We maintain the customers ordered in increasing order of their arrival times, so that if there are more than  $N$  customers in the system (where  $N$  is the number of agents) the first  $N$  elements represent the customers that are in service. Under these definitions  $\Xi^\Lambda(t)$  is clearly a Markov process.

We are now ready to prove proposition 1 which is more formally stated as follows: Fix  $\Lambda$  and assume  $N = R(1 + z)$  for some  $z > 0$ . Then, the process  $\Xi^\Lambda(t)$  admits a unique stationary distribution,  $\nu^\Lambda$ , which is also the limit distribution, that is

$$\Xi^\Lambda(t) \Rightarrow \Xi^\Lambda(\infty), \text{ as } t \rightarrow \infty,$$

where  $\Xi^\Lambda(\infty)$  has the distribution  $\nu^\Lambda$ , and the convergence holds regardless of the distribution of  $\Xi^\Lambda(0)$ . We actually prove a stronger result than mere stability under (S)-(C) with finite thresholds. That is, we prove that there exists a unique limit distribution even when we assume that all thresholds are set to  $\infty$ . The proof can then be easily adapted to cover finite thresholds. Since  $\Lambda$  is fixed we omit it from the notation throughout the proof. What we need to show to establish the statement of the proposition is that there exists  $\delta^0$  and  $K$  so that for all  $\xi \in \Xi$  with  $|\xi| > K$  we have that

$$E_\xi[|\Xi(|\xi|(1 + \delta^0))|] \leq 1/2|\xi|. \quad (20)$$

Assuming (20) holds one can follow the proof of Theorem 3.1 in Dai [10] to show positive Harris recurrence of  $\Xi(\cdot)$ . Positive Harris recurrence implies existence and uniqueness of a stationary distribution. The stationary and limit distributions agree by the obvious non-lattice structure of the time between consecutive visits to the origin (the 0 state of  $\mathbb{R}_+^\infty \times \mathbb{Z}_+$ ). The inequality (20) is established easily by adapting the fluid limits argument of [10] to our model. For the sake of completeness we give the detailed proof in the online appendix [17].  $\blacksquare$

We now turn to prove the results regarding the performance analysis and asymptotic optimality of (S)-(C). Here, the superscript  $\Lambda$  is used to denote parameters or processes related to the  $\Lambda^{th}$  system. We omit the superscript when dealing with parameters that do not scale with  $\Lambda$ .

In the rest of this appendix we use a different sample path construction than the one we used before and, in particular, one that relies on the Poisson nature of most of the processes involved. A detailed description of the construction is given in the online appendix [17]. Throughout the rest of the section we assume that the thresholds  $\eta_i^\Lambda$ ,  $i \leq \bar{k}$  satisfy  $\eta_i^\Lambda = \hat{\eta}_i \sqrt{\Lambda}$  for some constants  $\hat{\eta}_1 \geq \hat{\eta}_2 \dots \geq \hat{\eta}_{\bar{k}}$ .

**Proof of Theorem 1:** Since  $\bar{\Pi}(\Lambda) \geq \Pi^*(\Lambda) \geq \hat{\Pi}(\Lambda)$ , it suffices to prove that  $\hat{\Pi}(\Lambda) = \bar{\Pi}(\Lambda) - O(\sqrt{\Lambda})$ . Note that since (C) was proven to admit steady state we necessarily have that  $\lambda_i x_i(C) =$

$\mu_{cs}E[Z_{i,2}^\Lambda(\infty)]$  so that the profit function is given by

$$\hat{\Pi}(\Lambda) = \sum_{i=1}^{\bar{k}} r_i \mu_{cs} E[Z_{i,2}^\Lambda(\infty)] - cR(1 + \bar{z}) - d\Lambda E[W^\Lambda(\infty)].$$

It follows that in order to show the result it suffices to prove that

$$E[Z_{i,2}^\Lambda(\infty)] = \frac{\lambda_i q_i}{\mu_{cs}} - O(\sqrt{\Lambda}), \text{ for all } i \leq \bar{k}, \text{ and } E[W^\Lambda(\infty)] = O(1/\sqrt{\Lambda}). \quad (21)$$

Indeed, given that (21) holds the proof is completed since

$$\bar{\Pi}(\Lambda) - \hat{\Pi}(\Lambda) = \sum_{i=1}^{\bar{k}} r_i (\lambda_i q_i - \mu_{cs} E[Z_{i,2}^\Lambda(\infty)]) + d\Lambda E[W^\Lambda(\infty)] = O(\sqrt{\Lambda}) = o(\Lambda). \quad (22)$$

Expression (21) is rather intuitive. The left hand side says that in steady-state the expected number of type  $i$  customers that are in the system and are being cross-sold is equal to  $\lambda_i q_i / \mu_{cs}$ , which is what the deterministic relaxation predicts, minus a small (second order) correction term. Equivalently, the stochastic effects in the system result in small deviations around the deterministic solution. Similarly, the right hand side says that the expected waiting times experienced by callers are of order  $1/\sqrt{\Lambda}$ , which in turn are consistent with the second order correction terms just mentioned. The proof of (21) is rather technical. The complexity emanates from the fact that, while one may use simpler methods to prove convergence of the related processes (see, e.g., [3]), proving the convergence of the steady state variables is rather involved even if one can compute the steady state distribution explicitly - which is not the case in our setting. Following Gamarnik and Zeevi [14], we use an appropriate Lyapunov function constructed via a fluid model analysis together with some probabilistic bounds obtained using tools from Strong Approximations. The details of this argument is relegated to the technical (online) appendix [17].

**Proof of Proposition 2:** Item i. of the proposition was already proved within the proof of Theorem 1. The rest of the result is proved by contradiction. Assume first that ii. does not hold. That is that

$$\limsup_{\Lambda \rightarrow \infty} \frac{|N^*(\Lambda) - R(1 + \bar{z})|}{f(\Lambda)} > 0, \quad (23)$$

for some sequence  $f(\Lambda)$  with  $f(\Lambda)/\sqrt{\Lambda} \rightarrow \infty$  as  $\Lambda \rightarrow \infty$ . Note that  $\Pi^*(\Lambda) \leq \bar{\Pi}(N^*(\Lambda), \Lambda)$ , where the latter stands for the solution of the deterministic relaxation when the staffing level is fixed

to  $N^*(\Lambda)$ .  $\bar{\Pi}(N^*(\Lambda), \Lambda)$  is obtained by solving the corresponding fractional Knapsack problem. Assume for now that

$$\limsup_{\Lambda \rightarrow \infty} \frac{N^*(\Lambda) - R(1 + \bar{z})}{f(\Lambda)} > 0. \quad (24)$$

That is, the staffing is higher than the one suggested by (S) (the proof for the other case is essentially the same). Then, there exists a subsequence  $\{\Lambda^j\}_{j \geq 1}$ , such that for all  $j$ , the solution to the fractional knapsack problem is obtained by setting

$$x_i^{\Lambda^j} = q_i, \forall i \leq \bar{k}$$

and

$$\sum_{i=\bar{k}+1}^K \lambda_i x_i^{\Lambda^j} = \min \left\{ \mu_{cs}(N^*(\Lambda^j) - R(1 + \bar{z})), \sum_{i=\bar{k}+1}^K \lambda_i q_i \right\}$$

and we define the sum above to be identically 0 if  $\bar{k} = K$ . Then, we have that

$$\begin{aligned} \Pi^*(\Lambda^j) - \bar{\Pi}(\Lambda^j) &\leq \bar{\Pi}(N^*(\Lambda^j), \Lambda^j) - \bar{\Pi}(\Lambda^j) = -c(N^*(\Lambda^j) - R(1 + \bar{z})) + \sum_{i=\bar{k}+1}^K \lambda_i^j r_i x_i^{\Lambda^j} \\ &\leq -c(N^*(\Lambda^j) - R(1 + \bar{z})) + r_{\bar{k}+1} \mu_{cs}(N^*(\Lambda^j) - R(1 + \bar{z})) \\ &= (-c + r_{\bar{k}+1} \mu_{cs})(N^*(\Lambda) - R(1 + \bar{z})), \end{aligned} \quad (25)$$

where we set  $r_{\bar{k}+1} = 0$  if  $\bar{k} = K$ . Recalling that, by assumption,  $\mu_{cs} r_{\bar{k}+1} < c$ , and the definition of  $f(\Lambda)$ , we must have that  $\limsup_{\Lambda \rightarrow \infty} \frac{\Pi^*(\Lambda) - \bar{\Pi}(\Lambda)}{\sqrt{\Lambda}} = -\infty$ . But we have already shown that  $\Pi^*(\Lambda) = \bar{\Pi}(\Lambda) - O(\sqrt{\Lambda})$  leading to a contradiction.

The proof of iii. follows in essentially the same manner where we now assume, to reach a contradiction, that for some  $i$ ,  $x_i = \bar{x}_i + f(\Lambda)$ , with  $\frac{|f(\Lambda)|}{1/\sqrt{\Lambda}} \rightarrow \infty$  as  $\Lambda \rightarrow \infty$ . More specifically, if  $i \leq \bar{k}$  we assume that  $x_i^\Lambda = q_i + f(\Lambda)$  with  $f(\Lambda)$  negative for all  $\Lambda$ , and if  $i \geq \bar{k}$  we assume that  $x_i = f(\Lambda)$  with  $f(\Lambda)$  positive for all  $\Lambda$ . We then consider again the deterministic relaxation where, instead of fixing the staffing level, we fix  $x_i^\Lambda$ . The rest of the argument is now analogous to the proof of ii. ■

**Proof of Lemma 2:** Let  $\Pi(N^*(\Lambda), x^*(\Lambda))$  be the resulting profit in a system with arrival rate

$\Lambda$  and equipped with the optimal policy. Then, it is trivial that,

$$\Pi(N^*(\Lambda), x^*(\Lambda)) \leq -cR + \sum_{i=1}^K \lambda_i q_i (r_i - c/\mu_{cs})^+ - \Lambda dE[W^{\Lambda,*}] = \bar{\Pi}(\Lambda) - \Lambda dE[W^{\Lambda}(N^*(\Lambda), x^*(\Lambda))]. \quad (26)$$

Assume, by contradiction, that

$$\limsup_{\Lambda \rightarrow \infty} \sqrt{\Lambda} E[W^{\Lambda,*}] = \infty.$$

Let  $(N'(\Lambda), x'(\Lambda))$  the staffing level and fraction of cross-selling attempts resulting from (S)-(C) and recall that we have established in Lemma 1 that with thresholds satisfying (6), we have that

$$E[W^{\Lambda}] = O\left(\frac{1}{\sqrt{\Lambda}}\right), \text{ and } \bar{\Pi}(\Lambda) - \Pi^{\Lambda}(N'(\Lambda), x'(\Lambda)) = O(\sqrt{\Lambda}),$$

where  $W^{\Lambda}$  is the steady state waiting time under (S) – (C). In particular,

$$\limsup_{\Lambda \rightarrow \infty} \frac{\bar{\Pi}(\Lambda) - \Pi(N^*(\Lambda), x^*(\Lambda))}{\bar{\Pi}(\Lambda) - \Pi(N'(\Lambda), x'(\Lambda))} = \infty, \quad (27)$$

contradicting the optimality of  $(N^*(\Lambda), x^*(\Lambda))$ . ■