

Research Paper No. 2037

**A Structural Model of Sales-Force Compensation
Dynamics: Estimation and Field Implementation**

**Sanjog Misra
Harikesh Nair**

August 2009

RESEARCH PAPER SERIES

STANFORD
GRADUATE SCHOOL OF BUSINESS



<http://ssrn.com/abstract=1474462>

A Structural Model of Sales-Force Compensation Dynamics: Estimation and Field Implementation

Sanjog Misra* Harikesh Nair†

First version: June 2008
This version: August 2009‡

Abstract

We present an empirical framework to analyze real-world sales-force compensation schemes. The model is flexible enough to handle quotas and bonuses, output-based commission schemes, as well as “ratcheting” of compensation based on past performance, all of which are ubiquitous in actual contracts. The model explicitly incorporates the dynamics induced by these aspects in agent behavior. We apply the model to a rich dataset that comprises the complete details of sales and compensation plans for a set of 87 sales-people for a period of 3 years at a large contact-lens manufacturer in the US. We use the model to evaluate profit-improving, theoretically-preferred changes to the extant compensation scheme. These recommendations were then implemented at the focal firm. Agent behavior and output under the new compensation plan is found to change as predicted. The new plan resulted in a 9% improvement in overall revenues, which translates to about \$0.98 million incremental revenues per month, indicating the success of the field-implementation. The results bear out the face validity of dynamic agency theory for real-world compensation design. More generally, our results fit into a growing literature that illustrates that dynamic programming-based solutions, when combined with structural empirical specifications of behavior, can help significantly improve marketing decision-making, and firms’ profitability.

*Simon School of Business, Rochester University, Email: sanjog.misra@simon.rochester.edu;

†Graduate School of Business, Stanford University, Email: harikesh.nair@stanford.edu.

‡We thank Dan Akerberg, Lanier Benkard, Adam Copeland, Paul Ellickson, Liran Einav, Wes Hartmann, Gunter Hitsch, Phil Haile, Sunil Kumar, Ed Lazear, Philip Leslie, Kathryn Shaw, Seenu Srinivasan, John Van Reenan, and seminar participants at Berkeley, Kellogg, NYU, Rochester, Stanford, UC Davis, as well as the Marketing Science, Marketing Dynamics, NBER-IO, SICS, SITE, and UTD FORMS conferences, for their helpful feedback. Finally, we thank the management of the anonymous, focal firm in the paper for providing data, for innumerable interviews, and for their support, without which this research would not have been possible. We remain, however, responsible for all errors, if any.

1 Introduction

Personal selling via sales-forces is an important part of the economy. In the US, nearly 12% of the total workforce is employed in full-time sales occupations (Zoltners et al. 2001). In a review of sales-force practice, Albers and Mantrala (2008) note, “Dartnell’s 30th *Sales Force Compensation Survey: 1998–1999* reports the average company spends 10% and some industries spend as much as 40% of their total sales revenues on sales force costs. In total, the US economy is estimated to spend \$800 billion on sales forces, almost three times the amount spent on advertising in 2006 (Zoltners et al. 2008)”. The academic literature has recognized this practitioner interest, and the design of plans to compensate sales-force is now one of the most visible and successful applications of agency theory in real-world business settings (Mantrala, Sinha and Zoltners 1994). Surprisingly however, the richness of the theory (reviewed later), contrasts sharply with the sparsity of empirical work on the topic, stemming partly from the lack of detailed data on agent’s compensation and sales. The need for more empirical work is accentuated by the importance of accounting for several important features of real-world compensation schemes in evaluating and optimizing sales-force performance. Actual schemes in practice tend to be discrete and kinked, featuring quotas, bonuses and ceilings. In a survey of Fortune 500 firms, Joseph and Kalwani (1998) report that 95% of compensation schemes they survey had some combination of quotas and commissions, or both. These aspects complicate the analysis of contracts by generating dynamics in the actions of agents. A proper accounting of these dynamics then becomes critical to the evaluation and improvement of the sales-force scheme.

The goals of this paper are two fold. First, we present a framework that can help evaluate the dynamic effects of compensation contracts on sales-agent output. Our approach is to develop a dynamic structural model of agent behavior which we use, along with compensation and sales data, to estimate structural primitives underlying agent behavior. We discuss how the variation induced by the dynamics, in combination with rich data on actual compensation are useful in helping us learn about these primitives. Our second goal is to demonstrate how the model can be used to improve extant compensation schemes. We demonstrate how the structural model may be used to simulate agent behavior and firms’ outcomes to uncover potentially better, sales-force plans. Our framework can handle the rich variety in quotas, bonuses and

commissions schemes observed in practice, and can help decision-makers improve these plans. We present an application to the studying the sales-force contracts of a Fortune 500 contact lens manufacturer, where the recommendations based on the model were actually implemented. The recommendations involved changes to the nature and slope of output based incentives, and were implemented starting January 2009. Under the new plan, revenues to the firm increased by \$0.98 million per month across all agents (about \$12 million incremental per year, a 9% improvement overall), indicating the field-implementation was a success.¹ We interpret these results as strongly reinforcing the practical value of structural agency-theoretic models for compensation design in real-world settings. More generally, our results fit into a new literature that has illustrated the value of dynamic structural models for improving marketing decision-making via direct field interventions (e.g. Mantrala et al. 2006; Cho and Rust 2008).

The compensation plan in our data features a straight salary paid out irrespective of effort, as well as a marginal compensation on any sales generated above a “quota”, and below a “ceiling”. Such quotas are ubiquitous in sales-force compensation and have been justified in the theory literature as a trade-off between the optimal provision of incentives versus the cost of implementing more complicated schemes (Raju and Srinivasan 1996), or as optimal under specific assumptions on agent preferences and the distribution of demand (Oyer 2000). Ceilings or caps on compensation have been rationalized as a response to demand uncertainty on the part of the firm. To the extent that compensation should reflect effort, a ceiling has the advantage of hedging the firm against the payout of “windfall” compensation to agents for sales that were unrelated to effort.

While quotas are advantageous, they can also generate inefficiencies due the incentives of agents to time the allocation of effort. For instance, in a “salary + commission” scheme such as ours, sales-agents who achieve the quota required for earning the commission in the current compensation cycle may have a perverse incentive to postpone additional effort to the future. This enables the agent to use the sales generated from the postponed effort to attain the quota in the next compensation cycle. Indeed, in some settings, it is possible that such intertemporal reallocation of effort or “gaming” may negate the effort-inducing benefits from utilizing output-based compensation schemes. Similarly, ceilings have the obvious disadvantage of

¹This effect is likely to be a lower-bound on account of the recession in the US economy in 2009.

destroying the desirable convexity of the plan (e.g., Basu et al. 1985), thereby reducing the incentive to expend effort at high levels of output. A priori therefore, the elimination or enhancement of quotas or ceilings in a plan can either be beneficial or detrimental to the firm. As this is an empirical question, one of the questions we address is to empirically evaluate whether the profitability of the firm may be improved by changes to the compensation plan along these dimensions. A second question is whether eliminating quotas and ceilings altogether may be better. A third question is how these changes may be implemented taking into account organizational and culture-based constraints faced by the firm. Together, these aspects serve to illustrate the economic and managerial implications of the framework. We focus on these aspects specifically as they are features of the plan in our data, but the reader should note that the proposed framework is flexible enough to consider other changes to the compensation plan as well, including the introduction of bonuses and the provision of other nonlinear, output-dependent commission schemes.²

The main challenge in the analysis arises from the need to account for the dynamics in agents' actions induced by the shape of the compensation schedule. The source of the dynamics is the inherent nonlinearity of the plan. Quotas and ceilings generate curvature in the relationship between compensation and output. The curvature generates an incentive for effort-bunching: more is gained by the agent by expending a large effort in one month, than by spreading the same effort across many months. This in turn implies that a forward-looking perspective drives agent's effort allocation decisions. A second dynamic arises because of a common practice termed "*ratcheting*" whereby quotas for future periods are updated according to the agent's currently observed performance. Such ratcheting has been documented in several real world compensation schemes (e.g. Weitzman 1980; Leone, Misra and Zimmerman 2004), and is also a feature of the plan used by the firm in our empirical application. Ratcheting implies that the agent's current effort has an effect on his payoffs in future quarters, thereby making his effort allocation problem dynamic. A careful consideration of the dynamics then becomes essential to the estimation of the agent's preferences, and the simulation of his behavior under alternative compensation plans.

Estimation of the model is complicated by the fact that effort is unobserved. We

²Solving for the *optimal plan* is outside of the scope of the current analysis, and is an important, but methodologically challenging, direction for future research.

introduce a methodology that exploits the richness of our data, an informative structure, and recent advances in estimation methods to facilitate the identification of this latent construct. In particular, following the intuition in Copeland and Monnett (2009), we describe how intertemporal linkages helps identify effort from sales data in sales-force compensation settings. We model agents as maximizing intertemporal utility, conditional on the current compensation scheme, and their expectations about the process by which quotas would be updated based on their chosen actions. Our empirical approach is to estimate, in a first stage, the structural parameters involving the sales person’s utility function. We then simulate, in a second stage, his behavior given a changed compensation profile. The estimator for the 1st stage of our empirical strategy is based on the recent literature on 2-step estimation of dynamic decisions (Hotz and Miller 1993; Bajari, Benkard and Levin 2007). Our approach is to nonparametrically estimate agent-specific policy functions, and use these, along with the conditions for the optimality of the observed actions, to estimate the structural parameters. We discuss how an individual rationality constraint as well as the assumption of agent optimality identifies agent preferences. We use our estimates to generate the empirical distribution of agent preferences, which we use to simulate the behavior of the agent-pool under counterfactual compensation profiles.

A practical concern with the use of two step estimators has been the presence of unobserved serially correlated state variables which prevent consistent nonparametric estimation of first-stage policy functions and transitions. In particular, this ruled out models with unobserved heterogeneity (though see Arcidiacono & Miller 2008 for a recent approach that handles discrete unobserved heterogeneity). We are able to address this problem due to the availability of panel data of relatively long cross-section and duration for each agent, which facilitates estimation agent-by-agent. This enables a nonparametric accommodation of unobserved heterogeneity. Given the estimates from the first stage, we evaluate agent behavior and sales under the counterfactual by solving the agents’ dynamic programming problem numerically. We believe we are the first in the empirical literature to model the intertemporal problem facing sales-agents and to measure the dynamic effect of quotas and ratcheting in a real world setting.

Our model-free analysis of the data reveals evidence that the current plan may be inefficient. In particular, we find evidence for shirking by agents in the early part of the compensation cycle. The model predicts that elimination of quotas, and

reduction of the length of the quota cycle reduces this perverse incentive; this aspect is borne out in the realized sales from the new plan. We also find evidence that the extent of demand uncertainty may not be high enough to warrant the ceiling imposed on incentive compensation in the current plan. Indeed, in the new plan, ceilings are eliminated, and realized sales significantly exceed the caps from the old plan, as predicted by the model. Overall, our prediction from the model is that overall sales will rise from the elimination of quotas and ceilings, which is validated by the data from the new plan. The model also predicts that under the new plan, output variation within the months of the old compensation cycle will be eliminated, i.e. shirking in the early sales will reduce, and effort from the later part of the cycle will be reallocated to earlier months. The new data corroborates these predictions. Further, the differences in sales across months is not statistically significant under the new plan. Overall, these results strongly establish the out of sample validity and statistical power of the predictions from the proposed model.

Our paper adds into a small empirical literature that has explored the dynamic effects of incentive schemes. Despite the preponderance of nonlinear incentive schemes in practice, the empirical literature analyzing these, and the effect of quotas on sales-force effort in particular, has remained sparse. Part of the reason for the paucity of work has been the lack of availability of agent-level compensation and output data. The limited empirical work has primarily sought to provide descriptive evidence of the distortionary effect of payment schemes on outcomes (e.g. Healy 1985, in the context of executive compensation; Asch 1990, in the context of army-recruiters; and Courty and Marschke 1997, in the context of federal job-training programs). Oyer (1998) was the first to empirically document the timing effects of quotas, by providing evidence of jumps in firms' revenues at the end of quota-cycles that are unrelated to demand-side factors. On the theory side, it is well known (Holmstrom 1979; Lazear 1986) that nonlinear output-based contracts, in general, have the beneficial effect of inducing agents to exert effort, even when effort is unobservable by the firm. However, surprisingly little is known about the role of quotas in motivating agents effort.³ In the sales-force context there is a large literature that investigates the design and implementation of compensation plans that induce optimal levels of sales-force effort,

³An alternative motivation of output-based contracts is that it may help attract and retain the best sales-people (Lazear 1986; Godes 2003; Zenger and Lazarini 2004). This paper abstracts away from these issues since our data does not exhibit any significant turnover in the sales-force.

and examines the role of various factors on the nature and curvature of the optimal contract (see for e.g. Basu et al. 1985; Lal and Srinivasan 1993; Rao 1990). Most of this literature, however, has little to say about quotas (Coughlan 1993).

A related literature also seeks to empirically describe the effect of incentives, more broadly, on output (e.g. Chevalier and Ellison 1999; Lazear 2000; Hubbard 2003; Bandiera, Baransky and Rasul 2005; see Pendergast 1999 for a review). We complement this literature by detecting and *measuring* the dynamic inefficiencies associated with compensation schemes. The descriptive evidence on quotas are mixed. Using data from a different context, and a different compensation scheme, Steenburgh (2008) reports that agents facing quotas in a durable-goods company do not tend to reduce effort in response to lump-sum bonuses. In contrast, Larkin (2006) uses reduced form methods to document the distortionary effects of compensation schemes on the timing and pricing of transactions in technology-markets. The differences accentuate the need for more empirical work. Our paper is also related to the work of Ferrall and Shearer (1999), Paarsch and Shearer (2000), Lee and Zenios (2007), and Jiang and Palmatier (2009), who estimate static, structural models of agent behavior, while modeling the optimal contract choice by the firm. The closest paper to ours in spirit is Copeland and Monnett (2009) who estimate a dynamic model to analyze the effects of nonlinear incentives on agents' productivity in sorting checks. Our institutional context, personal selling by sales-force agents, adds several aspects that warrant a different model, analysis, and empirical strategy from Copeland and Monnet's context. Unlike their industry, demand uncertainty plays a key role in our setting; this generates a role for risk aversion, and a trade-off between risk and insurance in our contracts. Further, ratcheting, an important dynamic affecting agent effort in our setting, is not a feature of their compensation scheme. Ratcheting generates a dynamic across compensation periods, in addition to dynamics induced within the period by the nonlinearity.

The methods we develop here can also be used to analyze compensation issues in other business contexts. For example, a recent working paper by Chung, Steenburgh and Sudhir (2009) analyzes the role of bonuses in a durable good selling context. In contrast to our application, their plan has lump-sum bonuses and a progressive incentive scheme. As more limited panel data is available per agent, they demonstrate how the algorithm proposed by Arcidiacono and Miller (2009) may be used to pool across agents to control for unobserved heterogeneity. Similar to Steenburgh (2008),

they find their plan produces no distortions in effort.

Finally, our paper also adds to the theoretical literature on sales-force compensation by offering a computational framework in which to examine more realistic comparative dynamics that involve arbitrarily complex and dynamic compensation plans and effort policies of agents that respond to these dynamics. The rest of this paper is structured as follows: We begin with a description of our data and some stylized facts. We then introduce our model followed by the estimation methodology. We then discuss results and predictions for an improved plan. We then discuss results from the field implementation and then conclude.

2 Patterns in the Data and Stylized Facts

In this section, we start by presenting some stylized facts of our empirical application, and also provide model-free evidence for the effect of quotas on the timing of effort allocations by sales-agents in our data. We use the reduced form evidence and the stylized facts presented here to motivate our subsequent model formulation and empirical strategy.

2.1 Data and Compensation Scheme

Our data come from the direct selling arm of the sales-force division of a large contact lens manufacturer in the US with significant market-share in the focal category (we cannot reveal the name of the manufacturer due to confidentiality reasons). Contact lenses are primarily sold via prescriptions to consumers from certified physicians. Importantly, industry observers and casual empiricism suggests that there is little or no seasonality in the underlying demand for the product. The manufacturer employs 87 sales-agents in the U.S. to advertise and sell its product directly to each physician (also referred to as a “client”), who is the source of demand origination. The data consist of records of direct orders made from each doctor’s office via a online ordering system, and have the advantage of tracking the timing and origin of sales precisely. Agents are assigned their own, non-overlapping, geographic territories, and are paid according to a nonlinear period-dependent compensation schedule. We note in passing that prices play an insignificant role since the salesperson has no control over the pricing decision and because price levels remained fairly stable during the period for

which we have data.⁴ As noted before, the compensation schedule involves salaries, quotas and ceilings. Commissions are earned on any sales exceeding quota and below the ceiling. The salary is paid monthly, and the commission, if any, is paid out at the end of the quarter. The sales on which the output-based compensation is earned are reset every quarter. Additionally, the quota may be updated at end of every quarter depending on the agent’s performance (“ratcheting”). Our data includes the complete history of compensation profiles and payments for every sales-agent, and monthly sales at the client-level for each of these sales-agents for a period of about 3 years (38 months).

Quarterly, kinked compensation profiles of the sort in our data are typical of many real world compensation schemes. Consistent with the literature, our conversations with the management at the firm revealed that the primary motivation for quotas and commissions is to provide “high-powered” incentives to the sales-force for exerting effort in the absence of perfect monitoring. We also learned that the motivation for maintaining a “ceiling” on the compensation scheme is consistent with the “windfall” explanation mentioned in the introduction. The latter observation suggests that unanticipated shocks to demand are likely important in driving sales.

The firm in question has over 15,000 SKU-s (Stock Keeping Units) of the product. The product portfolio reflects the large diversity in patient profiles (e.g. age, incidence of astigmatism, nearsightedness, farsightedness or presbyopia, corneal characteristics, eye-power etc.), patient needs (e.g. daily, disposable, sports use, cleaning frequency) and contact lens characteristics (e.g. hydrogel, silicone-hydrogel, moistness, color etc.). The product portfolio of the firm is also characterized by significant new product introduction and line extensions reflecting the large investments in R&D and testing in the industry. New product introductions and line extensions reflect both new innovations as well as new usage regimens for patients uncovered by fresh trials and testing. The role of the sales-agent is primarily *informative*, by providing the doctor with updated information about new products available in the product-line, and by suggesting SKU-s that would best match the needs of the patient profiles currently faced by the doctor. While agents’ frequency of visiting doctors is monitored by the firm, the extent to which he “sells” the product once inside the doctor’s

⁴In other industries, agents may have control over prices (e.g. Bharadwaj 2002). In such situations, the compensation scheme may also provide incentives to agents to distort prices to “make quota”. See Larkin (2006), for empirical evidence from the enterprise resource software category.

office cannot be monitored or contracted upon. In addition, while visits can be tracked, whether a face-to-face interaction with a doctor occurs during a visit in within the agent’s control (e.g., an unmotivated agent may simply “punch in” with the receptionist, which counts as a visit, but is low on effort).⁵ In our application, we do not separately model these dimensions of sales-calls, and interpret all factors by which an agent shifts a doctors’ sales as effort.

2.2 The Timing of Effort

We start by checking whether dynamics are an important consideration for understanding agents’ behavior under this contract. We start by looking in the data to see whether there exists patterns consistent with agent’s shifting the allocation of sales within the compensation cycle in manners consistent with incentives. First, as Oyer (1998) pointed out, when incentives exist for agents to manipulate timing, output (i.e. sales) should look lumpy over the course of the sales-cycle. In particular, we expect to see spikes in output when agents are close to the end of the quarter (and most likely to be close to “making quota”). Figure 1 plots the sales achieved by month of quarter across sales-agents. Figure 1 reveal significant increase at the end of quarters suggesting that agents tend to increase effort as they reach closer to quota. In the absence of seasonality, this suggests the possibility of shirking early in the quarter.

In Figure 2, we present plots at the agent-level that suggest that agents also tend to reduce effort within the quarter. We plot patterns in sales (normalized by total sales across all months in the data) for four agents. The shaded regions in Figure 2 highlights quarters in which sales *fell* in the last month of the quarter, perhaps because the agent realized a very large negative shock to demand early in the quarter and reduced effort, or because he “made quota” early enough, and hence decided to postpone effort to the next sales-cycle.⁶ We now explore how these sales-patterns are related to how far the agent is from his quarterly quota. Figure 3 shows nonparametric estimates of the relationship between sales (y -axis) and the distance to quota (x -axis), computed across all the sales-people for the first two months of each quarter in the data. We define the distance to quota as $\frac{(\text{Cumulative Sales at beginning of month-quota})}{\text{quota}}$.

⁵The firm does not believe that sales-visits are the right measure of effort. Even though sales-calls are observed, the firm specifies compensation based on sales, not calls.

⁶One alternative explanation for these patterns is that the spikes reflect promotions or price changes offered by the firm. Our extensive interactions with the management at the firm revealed that prices were held fixed during the time-period of the data (in fact, prices are rarely changed), and no additional promotions were offered during this period.

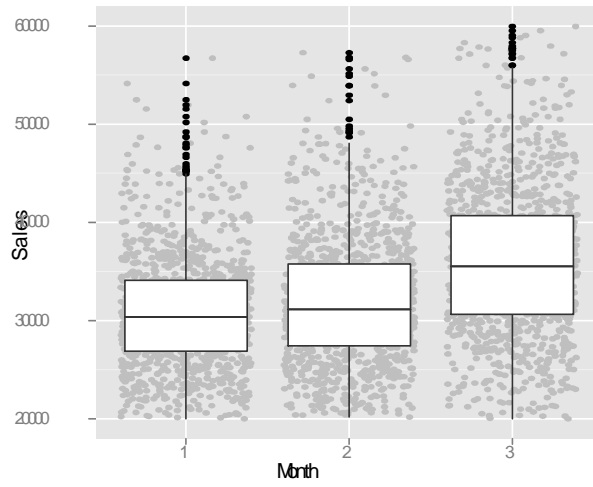
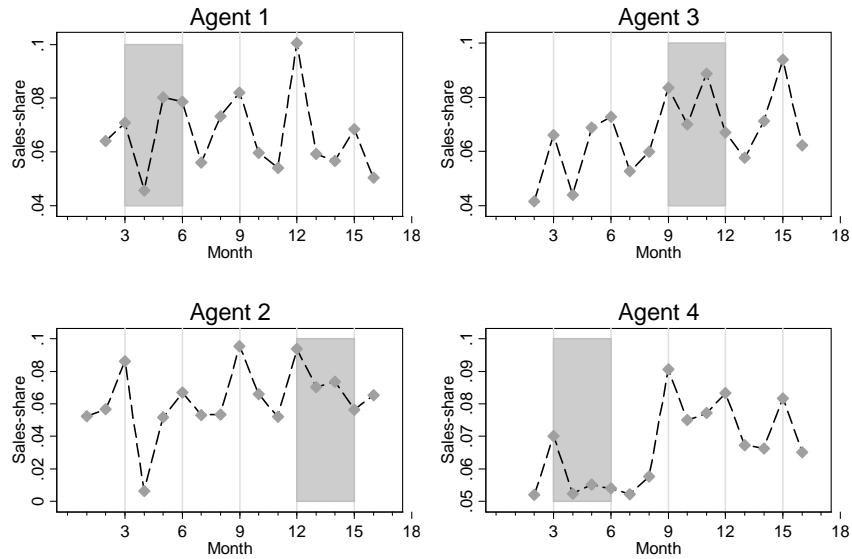


Figure 1: Sales per week by month of quarter.



Patterns suggesting agent 'gave up'

Figure 2: Agents reduce effort within quarters.

From Figure 3, we see that the distance to quota has a significant influence on the sales profile. Sales (proportional to effort) tend to increase as agents get closer to quota, suggesting increasing effort allocation, but fall once the agent reaches about 40% of the quota in the first 2 months, suggesting the agent anticipates he would “make the quota” by the end of the quarter. The decline in sales as the agent approaches quota is also consistent with the ratcheting incentive, whereby the agent reduces effort anticipating his quota may be increased in the next cycle, if he exceeds the ceiling this quarter. To further explore the effect of quotas, we present in Figure 4,

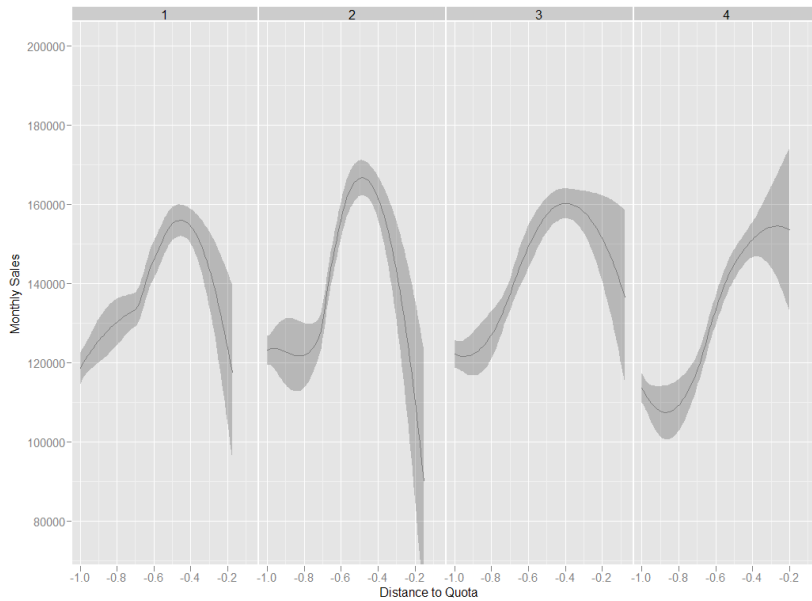


Figure 3: Sales vs distance to quota.

nonparametric plots of the % quota attained by the end of month $T - 1$ versus the % quota attained by the end of month T ($= 2, 3$), across all agents and quarters. Figure 4 suggests patterns that are consistent with intertemporal effort allocation due to quotas. In particular, when far away from quota in month $T - 1$ ($x \in 0.2, 0.4$), the profile is convex, suggesting a ramping up of effort. When the agent is close to quota in month $T - 1$ ($x \in 0.5, 0.8$), the profile is concave suggesting a reduction in the rate of effort allocation. Finally, figure 4 also shows that most agents do not achieve sales more than $1.4 \times$ quota, which is consistent with the effect of the ceiling (which was set to be $1.33 \times$ quota by the firm during the time-period of the data).

Figure 5 presents the analogous relationship, with plots for each agent in the data.

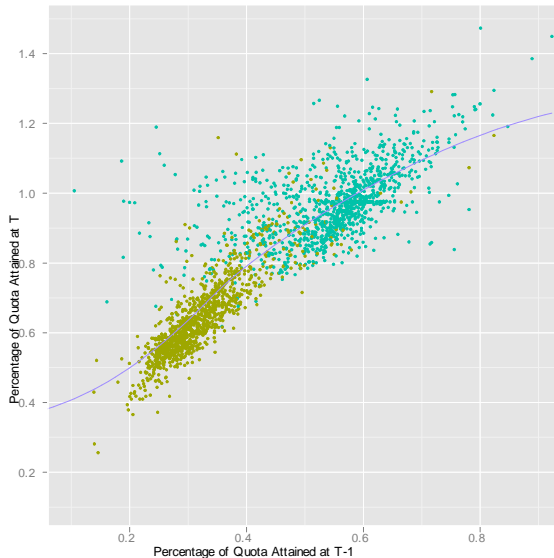


Figure 4: Concavity in Quota Attainment within Each Quarter.

Figure 5 shows that the concavity that we uncover is robust, and is not driven by pooling across agents.

Finally, we now consider whether these patterns are due to alternative phenomena unrelated to the effects of compensation schemes. The two leading explanations are a) demand side seasonality; and b) buyer side stockpiling. In the remainder of this section, we discuss how the institutional features of our setting, as well as the availability of some additional data enable us to rule out these explanations.

A priori, seasonality is not a compelling consideration due to the fact that the disease condition that the product treats is non-seasonal. Patient demand for the product tends to be flat over the year. Our extensive discussions with sales-agents as well as the management at the firm suggest that stockpiling by clients (i.e. doctors) is also not a relevant consideration in this category. First, as noted before, there are a large number of SKU-s available from the firm (about 15,000). The doctor is concerned about patient satisfaction and health, both of which are strongly linked to finding an exact match between the patient's needs and the right SKU from this large product set. Ex ante, the distribution of patient profiles, needs and usage characteristics arriving at his office for the coming month is uncertain. These considerations precludes stockpiling of SKU-s at the doctors office. The firm solves this supply-

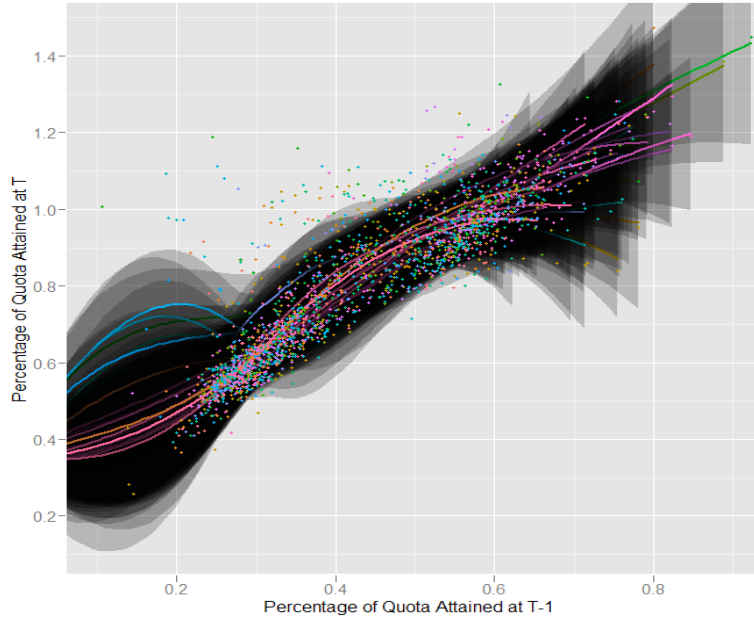


Figure 5: Concavity in Quota Attainment within Each Quarter by Agent.

chain problem by shipping the product directly to the consumer from its warehouse, upon receipt of an online order from the doctor made at his office during a patient appointment.

These aspects are also borne out in the data. Specifically, regressions (available on request) of sales on lagged sales indicate that lagged sales, as well as functions of ordering history are not significant in explaining doctor’s current orders. This is inconsistent with the stockpiling story. To check for seasonality, we exploit some limited additional information available to us on the number of sales-calls made by each agent at each client every month. This “call” information is available for the last 18 months of the data. The sales-call is not a decision variable for the agent as neither the number of calls nor the allocation of calls across clients is under the control of the agent (i.e., these are set exogenously by the firm). To test for seasonality, we use the behavior of doctors when exposed to zero calls. We find that the sales generated across months with no calls are not statistically significantly different from one another, implying no seasonality. Finally, a strong test that incentives drive these patterns is provided by the data from the new plan implemented at the firm. Under the new plan, incentives for sales-agents to time effort across months of the quarter

are eliminated. Hence, the monthly differences should be eliminated if only incentives, and not buyer-side seasonality or intertemporal substitution, are the source of the time variation. This is indeed the case: sales generation under the new plan is found to be flat across months (please see §6.1.1).

Taken together, these results suggest that seasonality and buyer intertemporal substitution are not significant considerations for these data. These features arise from the specifics of our empirical setting. We anticipate that both these aspects are likely to be important in other situations, for instance, those involving durable-good selling, where intertemporal substitution is well known to be significant, or in B2B situations where the buyers are large firms, for which quarterly financial deadlines are known to be an important source of buyer-side seasonality in orders (Oyer 1998). Taken together, the above model-free evidence also point to the existence of significant effects of the compensation scheme on agent’s intertemporal effort allocations in these data, and motivates the dynamics incorporated into the model of agent effort.

2.2.1 Discussion

Our above discussion highlights three facts regarding salesperson effort: (i) Salespeople are forward looking in that they allocate current effort in anticipation of future rewards; (ii) they act in response to their current quarter compensation environment by increasing and reducing effort relative to their quarter goals; and, (iii) salespeople take into account the impact of their current actions on subsequent changes in future firm compensation policies. These facts will play key roles in the development of our formal model of dynamic effort allocation. We discuss this next.

3 A Model of Dynamic Effort Allocation

We consider the intertemporal effort allocation of an agent facing a period-dependent, nonlinear compensation scheme. The compensation scheme involves a salary, α_t , paid in month t , as well as a commission on sales, β_t . The compensation scheme is period-dependent in the sense that it specifies that sales on which the commission is accrued is reset every N months. The compensation scheme is nonlinear in the sense that the commission β_t may depend discontinuously on the extent to which his total sales over the sales-cycle, Q_t , exceeds a quota, a_t , or falls below a ceiling b_t . The extent to which the ceiling is higher than the quota determines the range of sales over which

the agent is paid the marginal compensation. While our framework is general enough to accommodate compensation schemes where $\{\alpha_t, \beta_t, a_t, b_t\}$ change over time, our empirical application has the feature that the salary, α and the commission-rate, β are time-invariant, and that the ceiling b_t is a known deterministic function of the quota a_t . We develop the model in the context of this simpler compensation plan. The choice of the structure of the incentive scheme by the firm is determined by reasons outside of our model. Our approach will be to solve for the agent's effort policy taking the firm's compensation policy as given, and to use the model to simulate agent-effort for counterfactual compensation profiles. Let I_t denote the months since the beginning of the sales-cycle, and let q_t denote the agent's sales in month t . Further, let χ_t be an indicator for whether the agent stays with the firm. $\chi_t = 0$ indicates the agent has left the focal company and is pursuing his outside option.⁷ The total sales, Q_t , the current quota, a_t , the months since the beginning of the cycle I_t , and his employment status χ_t are the state variables for the agent's problem. We collect these in a vector $\mathbf{s}_t = \{Q_t, a_t, I_t, \chi_t\}$, and collect the observed parameters of his compensation scheme in a vector $\Psi = \{\alpha, \beta\}$.

3.1 Actions

At the beginning of each period, the agent observes his state, and chooses to exert effort e_t . Based on his effort, sales q_t are realized at the end of the period. We assume that the sales production function satisfies three conditions.

1. Current sales is a strictly increasing function of current effort.
2. Current sales are affected by the state variables only through their effect on the agent's effort.
3. Unobservable (to the agent) shocks to sales are additively separable from the effect of effort.

Condition 1 is a fairly innocuous restriction that more effort result in more sales. Monotonicity of the sales function in effort enables inversion of the effort policy function from observed sales data. Condition 2 implies that the quota, cumulative sales or months of the quarter do not have a direct effect on sales, over and above their effect

⁷We assume that once the agent leaves the firm, he cannot be hired back (i.e. $\chi_t = 0$ is an absorbing state).

on the agent’s effort. As is discussed in more detail below, this “exclusion” restriction facilitates nonparametric identification of effort from sales data. Condition 2 rules out reputation effects for the agent (the fact that an agent has achieved high sales in the quarter does not make him more likely to achieve higher sales today); and also rules out direct end-of-the-quarter effects on sales (we find support for these restrictions in our data). Condition 3 is a standard econometric assumption. Based on the above, we consider sales-functions of the form,

$$q_t = g(e_t; z, \mu) + \varepsilon_t \quad (1)$$

where, $g(\cdot)$ is the sales production function, such that $\frac{\partial g(e; \mu)}{\partial e} > 0$, μ is a vector of parameters indexing $g(\cdot)$; z is a vector of observed factors (such as the number and type of clients in an agent’s sales-territory) that affects his demand; and ε_t is a mean-zero agent and month specific shock to demand that is realized at the end of the period, which is unobserved by the agent at the time of making his effort decision. We assume that ε_t is distributed i.i.d. over agents and time-periods with distribution $\mathcal{G}_\varepsilon(\cdot)$, to the estimated from the data. ε_t serves as the econometric error term in our empirical model (we present our econometric assumptions in detail in §4.1). In our empirical work, we will consider specifications in which the production function $g(\cdot)$ is heterogeneous across agents. For now, we suppress the subscript “ i ” for agent for expositional clarity.

3.2 Per-period utility

The agents’ utility is derived from his compensation, which is determined by the incentive scheme. We write the agent’s monthly wealth from the firm as, $W_t = W(\mathbf{s}_t, e_t, \varepsilon_t; \mu, \Psi)$. We model his utility each month as derived from the wealth from the firm minus the cost of exerting effort. We denote the cost function as $C(e_t; d)$, where d is a parameter to be estimated. We assume that agents are risk-averse, and that conditional on $\chi_t = 1$, their per-period utility function is,

$$u_t = u(Q_t, a_t, I_t, \chi_t = 1) = E[W_t] - r \text{ var}[W_t] - C(e_t; d) \quad (2)$$

Here, r is a parameter indexing the agent’s risk aversion, and the expectation and variance of wealth is taken with respect to the demand shocks, ε_t . The specification in equation (2) is attractive since it can be regarded as a second order approximation

to an arbitrary utility function.⁸ We now discuss the transition of the state variables that generate the dynamics in the agent's effort allocation problem. The payoff from leaving the focal firm and pursuing the outside option is normalized to zero,

$$u_t = u(Q_t, a_t, I_t, \chi_t = 0) = 0 \quad (3)$$

3.3 State Transitions

There are two sources of dynamics in the model. The nonlinearity in the compensation scheme generates a dynamic into the agent's problem because reducing current effort increases the chance to cross, say, the quota threshold tomorrow. A second dynamic is introduced since the agent's current effort also affects the probability that his compensation structure is updated in the future. Hence, in allocating his effort each period, the agent also needs to take into account how current actions affect his expected future compensation structure. These aspects are embedded in the transitions of the state variables in the model. In the remainder of this section, we discuss these transitions. Subsequently, we present the value functions that encapsulate the optimal intertemporal decisions of the agent.

The first state variable, total sales, is augmented by the realized sales each month, except at the end of the quarter, when the agent begins with a fresh sales schedule, i.e.,

$$Q_{t+1} = \begin{cases} Q_t + q_t & \text{if } I_t < N \\ 0 & \text{if } I_t = N \end{cases} \quad (4)$$

We assume that the agent has rational expectations about the transition of his quota, a_t . We use the observed empirical data on the evolution the agent's quotas to obtain the transition density of quotas over time. We estimate the following transition function that relates the updated quota to the current quota, as well as the performance of the agent relative to that quota in the current quarter,

$$a_{t+1} = \begin{cases} a_t & \text{if } I_t < N \\ \sum_{k=1}^K \theta_k \Gamma(a_t, Q_t + q_t) + v_{t+1} & \text{if } I_t = N \end{cases} \quad (5)$$

In equation (5) above, we allow the new quota to depend flexibly on a_t and $Q_t + q_t$, via a K -order polynomial basis indexed by parameters, θ_k . We use this flexible polynomial to capture in a reduced-form way, the manager's policy for updating agents' quotas. The term v_{t+1} is an i.i.d. random variate which is unobserved by the agent

⁸In case of the standard linear compensation plan, exponential CARA utilities and normal errors this specification corresponds to an exact representation of the agent's certainty equivalent utility.

in month t . The distribution of v_{t+1} is denoted $\mathcal{G}_v(\cdot)$, and will be estimated from the data. Allowing for v_{t+1} in the transition equation enables us to introduce uncertainty into the agent's problem. In our empirical work, we extensively test different specifications for the ratcheting policy, and provide evidence that the associated errors v_{t+1} are not serially correlated in the specifications we use. Lack of persistence in v_{t+1} implies that all sources of time-dependence in the agent's quota updating have been captured, and that the remaining variation is white noise.⁹

The transition of the third state variable, months since the beginning of the quarter, is deterministic,

$$I_{t+1} = \begin{cases} I_t + 1 & \text{if } I_t < N \\ 1 & \text{if } I_t = N \end{cases} \quad (6)$$

Finally, the agent's employment status in $(t + 1)$, depends on whether he decides to leave the firm in period t . The employment state tomorrow is thus a control variable for the agent today, and is described below.

3.4 Optimal Actions

Given the above state-transitions, we can write the agent's problem as choosing effort to maximize the present-discounted value of utility each period, where future utilities are discounted by the factor, ρ . We collect all the parameters describing the agent's preferences and transitions in a vector $\Omega = \{\mu, d, r, \mathcal{G}_\varepsilon(\cdot), \mathcal{G}_v(\cdot), \theta_{k,k=1,\dots,K}\}$. In month $I_t < N$, the agent's present-discounted utility under the optimal effort policy can be represented by a value function that satisfies the following Bellman equation,

$$V(Q_t, a_t, I_t, \chi_t; \Omega, \Psi) = \max_{\chi_{t+1} \in (0,1), e > 0} \left\{ \begin{array}{l} u(Q_t, a_t, I_t, \chi_t, e; \Omega, \Psi) \\ + \rho \int_\varepsilon V(Q_{t+1} = Q(Q_t, q(\varepsilon_t, e)), a_{t+1} = a_t, I_t + 1, \chi_{t+1}; \Omega, \Psi) f(\varepsilon_t) d\varepsilon_t \end{array} \right\} \quad (7)$$

The value in period $I_t + 1$ is stochastic from period I_t 's perspective because the effort in period I_t is decided prior to the realization of ε_t , which introduces uncertainty into the cumulative sales attainable next period. Hence, the Bellman equation involves an expectation of the $(I_t + 1)$ -period value function against the distribution of ε_t ,

⁹We also reject correlation of v_{t+1} across agents, as well as correlation of v_{t+1} with the demand shocks (ε_t) across agents. This rules out a story where subjective quota updating is used as a mechanism to filter out common shocks.

evaluated at the states tomorrow. Similarly, the Bellman equation determining effort in the last period of the sales-cycle is,

$$V(Q_t, a_t, N, \chi_t; \Omega, \Psi) = \max_{\chi_{t+1} \in (0,1), e > 0} \left\{ \begin{array}{l} u(Q_t, a_t, N, \chi_t, e; \Omega, \Psi) \\ + \rho \int_v \int_\varepsilon V(Q_{t+1} = 0, a_{t+1} = a(Q_t, q(\varepsilon_t, e), a_t, v_{t+1}), 1, \chi_{t+1}; \Omega, \Psi) \\ \times f(\varepsilon_t) \phi(v_{t+1}) d\varepsilon_t dv_{t+1} \end{array} \right\} \quad (8)$$

At the end of the sales-cycle, the cumulative sales is reset and the quota is updated. The value in the beginning of the next cycle is again stochastic from the current perspective on account of the uncertainty introduced into the ratcheted future quota by the demand shock, ε_t , and the quota-shock, v_{t+1} . Hence, the Bellman equation in (8) involves an expectation of the 1st period value function against the distribution of both ε_t and v_{t+1} .

Conditional on staying with the firm, the optimal effort in period t , $e_t = e(\mathbf{s}_t; \Omega, \Psi)$ maximizes the value function,

$$e(\mathbf{s}_t; \Omega, \Psi) = \arg \max_{e > 0} \{V(\mathbf{s}_t; \Omega, \Psi)\} \quad (9)$$

The agent stays with the firm if the value from employment is positive, i.e.,

$$\chi_{t+1} = 1 \text{ if } \max_{e > 0} \{V(\mathbf{s}_t; \Omega, \Psi)\} \geq 0$$

Given the structure of the agent's payoffs and transitions, it is not possible to solve for the value function analytically. We solve for the optimal effort policy numerically via modified policy iteration. The state-space for the problem is discrete-continuous, of dimension $\mathbb{R}^2 \times (N + 1)$. The two continuous dimensions (Q_t and a_t) are discretized, and the value function is approximated over this grid for each discrete value of N and employment status. One iteration of the solution took 120 seconds on a standard Pentium PC. Further computational details of the algorithm are provided in Appendix A. We now present the technique for the estimation of the model parameters.

4 Empirical Strategy and Estimation

Our empirical strategy is motivated by the intended use of the model, which is to obtain a relative evaluation of the outcomes for the firm under a changed compensation scheme. This requires a method to simulate the outcomes for the firm under

new compensation schemes. Consider a new compensation plan $\wp(q(e); \Psi)$, where Ψ indexes the parameters governing the features of the new plan (e.g. a revised salary, bonus, commission rate, quota etc.), and $q(\cdot)$ is sales.¹⁰ The firm's present discounted payoffs under $\wp(q(e); \Psi)$ are,

$$\Pi_\wp = \int \int \sum_{\tau=0}^{\infty} \beta^\tau [q(e_\wp) - \wp(q(e_\wp))] d\mathcal{F}(\mu, r, d) d\mathcal{G}_\varepsilon(\varepsilon_\tau) \quad (10)$$

where (e_\wp) is the effort policy expended by the agent when faced with compensation policy $\wp(q(e))$,

$$e_\wp = \arg \max_{e>0} V(s; e | \{\mu, r, d\}, \wp(\cdot))$$

In equation (10), $\mathcal{F}(\mu, r, d)$ is the joint CDF across agents in the firm of demand parameters, risk aversion and the cost of effort. Our approach will be to use the model to simulate effort and sales under the counterfactual plans, conditioning on estimates of $\mathcal{F}(\mu, r, d)$ and $\mathcal{G}_\varepsilon(\varepsilon_\tau)$.¹¹ A comparison of current policy quantities $\{\Pi^*, q^*, e^*\}$ to the counterfactual then facilitates a relative evaluation of the current plan to other potentially, better possibilities. The key object of econometric inference is thus the joint distribution of preferences, $\mathcal{F}(\mu, r, d)$ and of demand uncertainty, $\mathcal{G}_\varepsilon(\varepsilon_\tau)$. In the section below, we discuss a methodology that delivers estimates of these distributions.

Our discussion below comprises two steps. In step 1, we discuss how we use the observed data on sales and compensation plans across agents to estimate the parameters of the agents' preferences, as well as the functions linking sales to effort. In step 2, we discuss how we use these parameters, along with our dynamic programming (henceforth DP) solution to simulate the agent's actions under counterfactual compensation profiles. In the remainder of this section, we first discuss our econometric assumptions, and then present details on the specific compensation scheme in our

¹⁰Implicitly, Ψ can be a function of the agent's characteristics, $\Psi \equiv \Psi(\mu, r, d, \mathcal{G}_\varepsilon(\cdot))$. For example, a counterfactual scheme could be characterized by a fixed salary and a commission specific to each agent. In this contract, the optimal salary and commission rate would be a function of the agent's preferences. We suppress the dependence of Ψ on these features for notational simplicity.

¹¹Implicitly, in equation (10), we assume that the distribution of demand shocks, $\mathcal{G}_\varepsilon(\cdot)$ stays the same under the counterfactual. In equation (10), we do not integrate against the ratcheting shocks $\mathcal{G}_v(\cdot)$, because all the counterfactual contracts we consider involve no ratcheting. Consideration of counterfactual contracts that involve ratcheting would require a model for agents' belief formation about quota updating under the new compensation profile, which is outside of the scope of the current analysis. Future research could consider solving for the *optimal* quota updating policy, under the assumption that agents' beliefs regarding ratcheting are formed rationally. See Nair (2007) for one possible approach to solving for beliefs in this fashion applied to durable good pricing.

data. Subsequently, we describe the procedure for estimation of the parameters of the model.

4.1 Econometric Assumptions

The econometric assumptions on the model are motivated by the nature of the data, as well as the intended procedure for estimation. The observed variation to be explained by the model is the correlation of sales across months with the distance to quotas, the changes in sales when quotas change, as well as the variation of sales across agents, which are a function of the agents' effort. The computational challenge in estimation derives from the fact that the model implies that each agent's effort, and consequently, their sales, are solutions to a dynamic problem that cannot be solved analytically.

One approach to estimation would be to nest the numerical solution of the associated DP into the estimation procedure. This would be significantly numerically intensive since the DP has to be repeatedly solved for each guess of the parameter vector. Instead, our estimation method builds on recently developed methods for two-stage estimation of dynamic models (e.g. Hotz and Miller 1993; Bajari, Benkard and Levin 2007, henceforth BBL), which obviates the need to solve the DP repeatedly. Under this approach, agents' policy functions - i.e., his optimal actions expressed as a function of his state - as well as the transition densities of the state variables are estimated nonparametrically in a first-stage; and subsequently, the parameters of the underlying model are estimated from the conditions for optimality of the chosen actions in the data. We face two difficulties in adapting this approach to our context. First, the relevant action - effort - is unobserved to the econometrician, and has to be inferred from the observed sales. This implies that we need a way to translate the sales policy function to an "effort policy function". Second, unobserved agent heterogeneity is likely to be significant in this context, since we expect agents to vary significantly in their productivity. The dependence of sales on quotas induced by the compensation scheme generates a form of state dependence in sales over time, which in the absence of adequate controls for agent heterogeneity generates well-known biases in the estimates of the effort policy. However, handling unobserved heterogeneity in the context of 2-step Hotz-Miller type estimators has been difficult to date (there has been recent progress on this topic; please see Arcidiacono and Miller 2008).

We address both issues in our proposed method. To handle the first issue, we make

a parametric assumption about the sales-production function. We discuss below why a nonparametric solution is not possible. We are able to handle the second issue due to the availability of sales-information at the agent-level of unusually large cross-section and duration, which enables us to estimate agent-specific policy functions, and to accommodate nonparametrically the heterogeneity across agents. We discuss the specific assumptions in more detail below.

4.1.1 Preliminaries

The model of agent optimization presented in §3 implies that the optimal effort each period is a function of only the current state \mathbf{s}_t . To implement a two step method, we thus need to estimate nonparametrically in a first-stage, the effort policy function, $e_t = \hat{e}(\mathbf{s}_t)$. The effort policy function is obtained parametrically from the sales-policy function. To see the need for a parametric assumption, recall from §3 that we consider sales-production functions of the form,

$$q_t = g(e_t(\mathbf{s}_t), z) + \varepsilon_t$$

For clarity, we suppress the variable z , as the argument below holds for each value of z . Let $f(\mathbf{s}_t) \equiv g(e_t(\mathbf{s}_t))$.

Remark 1 *If at least two observations on q are available for a given value of \mathbf{s} , the density of $f(\mathbf{s})$ and ε are separately nonparametrically identified (Li and Vuong 1998).*

Remark 2 *Given the density of $f(\mathbf{s})$, only either $g(\mathbf{s})$ or $e(\mathbf{s})$ can be estimated nonparametrically.*

Remark 2 underscores the need for a parametric assumption on the relationship between sales and effort. One option to relax this would be to obtain direct observations on agent's effort, via say, survey data, or monitoring. This of course, changes the character of the principal-agent problem between the agent and the firm. Unobservability of agent effort is the crux of the moral hazard problem in designing compensation schemes. Hence, we view this parameterization as unavoidable in empirical models of sales-force compensation.

We now discuss how we use this assumption, along with the sales data to estimate the sale-production function. For each agent in the data, we observe sales at each

of J clients, for a period of T months. In our empirical application T is 38 (i.e., about 3 years), and J is of the order of 60-300 for each agent. The client data adds cross-sectional variation to agent-level sales which aids estimation. To reflect this aspect of the data, we add the subscript j for *client* from this point onward. In light of remark 2 we assume that the production function at each client j is linear in effort,

$$q_{jt} = h_j + e_t + \varepsilon_{jt} \quad (11)$$

$$= h_j(z_j) + e(\mathbf{s}_t) + \varepsilon_{jt} \quad (12)$$

The linear specification is easy to interpret: h_j can be interpreted as the agent’s time-invariant intrinsic “ability” to sell to client j , which is shifted by client characteristics z_j . We now let $h_j \equiv \mu'z_j$, and let $\hat{e}(\mathbf{s}_t) = \lambda'\boldsymbol{\vartheta}(\mathbf{s}_t)$, where γ is a $R \times 1$ vector of parameters indexing a flexible polynomial basis approximation to the monthly effort policy function. Then, the effort policy function satisfies,

$$q_{jt} = \mu'z_j + \lambda'\boldsymbol{\vartheta}(\mathbf{s}_t) + \varepsilon_{jt} \quad (13)$$

We assume that ε_{jt} is distributed i.i.d. across clients. We can then obtain the demand parameters and the effort policy function parameters from the following minimization routine,

$$\min_{\mu, \lambda} \|q_{jt} - (\mu'z_j + \lambda'\boldsymbol{\vartheta}(\mathbf{s}_t))\| \quad (14)$$

As a by product, we also obtain the effort policy function for the month t for each client as,

$$\hat{e}_t = \hat{\lambda}'\boldsymbol{\vartheta}(\mathbf{s}_t) \quad (15)$$

and the time-specific error distribution,

$$\hat{\varepsilon}_t = \sum_j \left(q_{jt} - \left(\hat{\mu}'z_j + \hat{\lambda}'\boldsymbol{\vartheta}(\mathbf{s}_t) \right) \right) \quad (16)$$

which is then used to estimate the empirical distribution of ε_t for each agent.¹² This distribution is an input to solving the dynamic programming problem associated with solution of the model for each agent. We sample with replacement from the estimated empirical distribution for this purpose.

¹²Alternatively, one could assume a parametric density for ε and use maximum likelihood methods. The advantage of our nonparametric approach is that we avoid the possibility of extreme draws inherent in parametric densities and the pitfalls that go along with such draws.

Finally, at the end of this step, we can recover the predicted overall sales for the agent which determines the agent’s overall compensation. Summing equation (13) across clients, the overall sales in month t is,

$$q_t = \sum_j^J q_{jt} = h + J e_t + \varepsilon_t \quad (17)$$

where, $h = \sum_{j=1}^J \hat{\mu}' z_j$, and $\varepsilon_t = \sum_{j=1}^J \hat{\varepsilon}_{jt}$. The total effort expended by a sales-agent across all clients in period t is thus $J e_t$, which affects per-period payoffs through equation (2)

Intuition for estimation of effort: Intuitively, we can think of identification of the effort policy by casting the estimator in equation (13) in two steps,

- Step 1: Estimate time-period fixed effects ϖ_t as, $q_{jt} = \mu' \mathbf{z}_j + \varpi_t + \varepsilon_{jt}$
- Step 2: Project ϖ_t on a flexible function of the state variables as $\varpi_t = \lambda' \boldsymbol{\vartheta}(\mathbf{s}_t)$

The client-level data facilitates the estimation of time-period specific fixed effects in Step 1. Equation (13) combines steps 1 & 2 into one procedure. We discuss the identification of the model in further detail below.

4.2 Compensation scheme

We now discuss the specifics of the compensation scheme in our dataset, and derive the expression for the monthly expected wealth for the agent given the above economic assumptions. The agent’s payout under the plan is determined based on his quarter-specific performance. Thus, $N = 3$, and cumulative sales, which affect the payout, are reset at the end of each quarter. The monthly salary α is paid out to the agent irrespective of his sales. If his current cumulative sales are above quota, the agent receives a percentage of a fixed amount β as commission. The percentage is determined as the proportion of sales above a_t , and below a maximum ceiling of b_t , that the agent achieves in the quarter. Beyond b_t , the agent receives no commission. For the firm in our empirical application, $\beta = \$5,000$, and the ceiling was always set 33% above the quota, i.e., $b_t = \frac{4}{3} a_t$. Figure 6 depicts the compensation scheme. We can write the agent’s wealth, $W(\mathbf{s}_t, e_t, \varepsilon_t; \mu, \Psi)$ in equation (2) as,

$$\begin{aligned}
W(\mathbf{s}_t, e_t, \varepsilon_t; \mu, \Psi) &= \alpha + \beta \left[\frac{(Q_t + q_t - a_t)}{b_t - a_t} \mathbf{I}(a_t \leq Q_t + q_t \leq b_t) + \mathbf{I}(Q_t + q_t > b_t) \right] \mathbf{I}(I_t = N) \\
&= \alpha + \beta \left[3 \frac{(Q_t + q_t - a_t)}{a_t} \mathbf{I}(a_t \leq Q_t + q_t \leq b_t) + \mathbf{I}(Q_t + q_t > b_t) \right] \mathbf{I}(I_t = N) \\
&= \alpha + \beta \min \left\{ \frac{3(Q_t + q_t - a_t)}{a_t}, 1 \right\} \mathbf{I}(Q_t + q_t > a_t) \mathbf{I}(I_t = N) \quad (18)
\end{aligned}$$

Thus, at the end of each sales-cycle, the agent receives the salary α , as well as a incentive component, $\beta \times \min \left\{ \frac{3(Q_t + q_t - a_t)}{a_t}, 1 \right\}$, on any sales in excess of quota. If it is not the end of the quarter, $\mathbf{I}(I_t = N) = 0$, and only the salary is received. Finally,

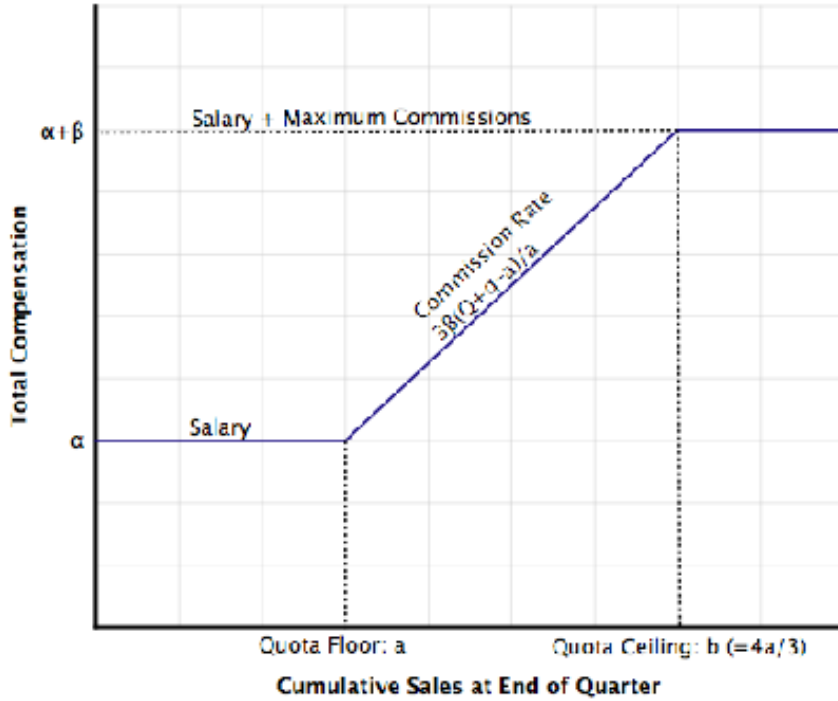


Figure 6: Compensation Scheme.

assume that the cost function in (2), $C(e)$, is quadratic in effort, i.e. $C(e_t) = \frac{de_t^2}{2}$, where d is a parameter to be estimated.

4.3 Estimation procedure

We now present the steps for estimation of the model parameters. The estimation consists of two steps, the first for a set of “auxiliary” parameters, and the second for

a set of “dynamic parameters.” We discuss these in sequence below.

4.3.1 Step 1: Nonparametric estimation of policy function and state transitions

The goals of the first step are two-fold. First, we estimate the demand parameters μ , as well as the distribution of demand shocks $\mathcal{G}_\varepsilon(\varepsilon_\tau)$ for each agent. Second, we estimate an effort policy function, as well as transitions of the state variables for each agent. We use both set of objects to estimate $\mathcal{F}(\mu, r, d)$ in step 2.

The effort policy function is related to observed sales via equation (13). The demand parameters and the demand shock distribution are obtained as by-products of estimating equation (13). We estimate the effort policy agent-by-agent. For each agent, data are pooled across the agent’s clients, and equation (13) estimated via least squares. An advantage of this approach is that we are able to handle heterogeneity across agents nonparametrically.

The next step is to estimate the parameters $(\theta_k, \mathcal{G}_v(\cdot))$ describing the transition of the agent’s quotas in equation (5). This is a series estimator which we estimate via nonlinear least squares. Since quotas vary only at the quarter-level, we do not estimate the quota transitions agent-by-agent. Instead, we pool the data across agents to estimate the quota transition function allowing for agent fixed-effects. The distribution of ratcheting shocks, $\mathcal{G}_v(\cdot)$, are estimated nonparametrically from the residuals from this regression.

The law of motion of the other state variables (month-of-the-quarter) does not have to be estimated since it does not involve any unknown parameters. This concludes step 1. Since we have estimated μ agent by agent, we can construct its marginal CDF ($\mathcal{F}(\mu)$) using a simple estimator,

$$\widehat{\mathcal{F}}(\mu) = \frac{1}{N} \sum_{i=1}^N \mathbf{I}(\mu_i \leq \mu). \quad (19)$$

The only remaining object to be estimated is the conditional distribution of the risk aversion, r , and the cost parameter, d , $\mathcal{F}(r, d|\mu)$. Step 2 below delivers estimates of $\mathcal{F}(r, d|\mu)$.

4.3.2 Step 2: Estimation of $\mathcal{F}(r, d|\mu)$

We estimate the “dynamic” parameters r and d using the methods proposed in BBL for the case of continuous controls. The BBL estimator is a minimum distance esti-

mator that finds parameters which minimize a set of moment inequality conditions. We propose to estimate the parameters by imposing two moment conditions that arise naturally in the class of principal-agent problems. In particular, let s_0 be an initial state for an agent, (r^*, d^*) be the true parameters, and e^* the optimal effort policy at the true parameters. Then, (r^*, d^*) must satisfy,

1. Individual Rationality (IR): $V(s_0; e^*, r^*, d^*) \geq 0$
2. Incentive Compatibility (IC): $V(s_0; e^*, r^*, d^*) \geq V(s_0; e', r^*, d^*)$

where $V(s_0; e^*, r^*, d^*)$ is the value function corresponding to the optimal policy e^* , and $V(s_0; e', r^*, d^*)$ is the present discounted utility corresponding to any other feasible policy, $e' \neq e^*$. The IR constraint says that the agent should at least be as better off working with the firm, as leaving the firm and pursuing his outside option. The IC constraint says that the agent should obtain higher utility in present discounted terms under the optimal effort policy, compared to any other feasible effort policy. Following BBL, we propose to estimate r^*, d^* by finding the set of parameters that minimize violations of these conditions over a random sample of the state space. In what follows, we assume that the optimal policy function $e^* = e^*(s_0)$ has already been estimated in step 1, and is available to the econometrician. Begin by defining the following quantities,

$$Z(s_0; e^*) = [\mathbb{E}(W) \quad \mathbb{V}(W) \quad \mathbb{C}(e)] \quad (20)$$

$$\theta = [1 \quad r \quad d] \quad (21)$$

where θ are the “dynamic” parameters to be estimated, s_0 is an initial state, e^* is the estimated optimal effort policy function and $Z(s_0; e^*)$ has components,

$$\mathbb{E}(W) = E_{e^*|s_0} \sum_{t=0}^{\infty} \beta^t E_{\varepsilon} [W(s, e^*(s))] \quad (22)$$

$$\mathbb{V}(W) = E_{e^*|s_0} \sum_{t=0}^{\infty} \beta^t E_{\varepsilon} [W(s, e^*(s))^2 - E_{\varepsilon} [W(s, e^*(s))]^2]$$

$$\mathbb{C}(e) = E_{e^*|s_0} \frac{1}{2} \sum_{t=0}^{\infty} \beta^t e^*(s)^2 \quad (23)$$

The value function based on the optimal effort policy can then be expressed as,

$$V(s_0; e^*, \theta) = Z(s_0; e^*)' \theta \quad (24)$$

Similarly, for any alternative policy function ($e' \neq e^*$), the perturbed value function is,

$$V(s_0; e', \theta) = Z(s_0; e')' \theta \quad (25)$$

Define the following two moment conditions,

$$\begin{aligned} g_1(s_0; \theta) &= \min(V(s_0; e^*, \theta), 0) \\ g_2(s_0, e'; \theta) &= \min(V(s_0; e^*, \theta) - V(s_0; e', \theta), 0) \end{aligned} \quad (26)$$

and let $g(s_0, e'; \theta) = [g_1(s_0; \theta) \quad g_2(s_0; \theta)]'$.

Let $H(\cdot)$ be a sampling distribution over states s_0 and alternative feasible policies e' . Define an objective function,

$$Q(\theta) = \int [g(s_0, e'; \theta)]' \Lambda [g(s_0, e'; \theta)] dH(s_0, e') \quad (27)$$

where, Λ is a 2×2 weighting matrix. Clearly, the true parameter vector ($\theta = \theta^*$) must satisfy,

$$Q(\theta^*) = \min_{\theta} (Q(\theta)) = 0 \quad (28)$$

Following BBL, we estimate θ^* by minimizing the sample analog of $Q(\theta)$. The function $Q(\theta)$ is obtained by averaging its evaluations over NR i.i.d. draws of s_0 from a uniform distribution over the observed support of states for the agent. At each s_0 , we generate alternative feasible policies by adding a normal error term to the estimated optimal effort policy. Using these, we forward simulate the terms in equation (22) to evaluate the moments at each guess of the parameter vector. The linearity of the value functions in θ imply that we can pre-compute $Z(s_0; e^*)$ and $Z(s_0; e')$ prior to parameter search, reducing computational time. In principle, an “optimal” Λ that weights each of the moment conditions based on their informativeness about θ would give the most efficient estimates. However, the econometric theory for the optimal Λ for inequality estimators of this sort are still to be developed. Hence, in practice, we set Λ equal to the identity matrix. This yields consistent but potentially inefficient estimates. Further computational details of our estimation procedure are presented in Appendix (A).

We perform estimation agent by agent. The main computational burden arises from forward-simulating value functions and implementing the nonlinear search separately for each agent (i.e. we solve 87 separate minimization problems). For each, we

obtain point estimates of $r, d|\mu$. We use these to construct a nonparametric estimate of the CDF across agents, $\mathcal{F}(r, d|\mu)$ as in the earlier section.

In general, the approach above yields point estimates of the parameters. Point estimation implies that the optimizer finds no other value of θ other than θ^* for which $Q(\theta) = 0$. A critical determinant to the point identification of the parameters is $H(\cdot)$. In particular, $H(\cdot)$ has to have large enough support over the states and alternative feasible effort policies to yield identification. This in turn requires that we a) pick the alternative feasible policies “intelligently”, such that they are informative of θ ; and b) more importantly, the econometrician has access to sufficient data (i.e. state points), on which nonparametric estimates of the optimal policy are available, and from which s_0 -s can be sampled. In application, we found that perturbations of the effort policy that were too far away from the effort policy were uninformative of the parameter vector. We use “small” perturbations (see Appendix (A) for precise details), which combined with the richness of our data, yield point identification of the parameters in our context for all the agents in the data.

4.4 Discussion: Identification

We now provide a more detailed discussion of identification in our model. In particular, we discuss how intertemporal linkages in observed sales identifies an agent’s unobserved effort allocation over time. The first concern is that effort has to be inferred from sales. In particular, looking at equation (11), we see that sales is explained by two unobservables, the first, effort, and the second, client-specific demand shocks. How can the data sort between the effects of either? The key identifying assumptions are,

1. Effort is a deterministic function of only the state variables.
2. Effort is not client specific - i.e., the agent allocates the same effort to each client in a given month.

We believe the first assumption is valid since we believe we have captured the key relevant state variables generating the intertemporal variation in agent effort. Further, after including a rich-enough polynomial in the state variables in equation (11), we can reject serial correlation in the residuals, ε_{jt} (i.e. the remaining variation in sales is only white noise). Assumption 1 is also consistent with our dynamic

programming model which generates a deterministic policy by construction. We believe the second assumption is reasonable. In separate analysis (not reported), we use limited data on the number of sales calls made by agents to each of the clients to check the validity of this assumption. In regressions of sales on calls, we find that the marginal effect of calls is not statistically significantly different across client-types, suggesting that effort more broadly, is not being tailored to each individual client.

Given these two assumptions, effort is identified by the joint distribution over time of the agent's current sales, and the extent to which cumulative sales are below or above the quota and the ceiling. To see this, recall that the optimal policy implies that the agent expends high effort when he is close to the quota, irrespective of month. The agent expends low effort when he has either crossed the ceiling in a given quarter, or when he is very far away from the quota in an early month. Under the former situation, the marginal benefit of an additional unit of effort is higher when expended in the next quarter; the same is true under the latter, since he has very little chance of reaching the quota in the current quarter. The model assumes that sales are strictly increasing in effort. Hence, if we see an agent achieve high sales across clients when he is close to the quota we conclude that effort is high. If we see low sales early on in the quarter, and when the quarter's sales have crossed the ceiling, we conclude that effort is low. Our identification argument is based essentially on the fact that variation in effort over time is related to variation in the distance to quota over time, and is similar to the identification of productivity shocks in the production economics literature (see e.g. Olley and Pakes 1996; Akerberg, Caves and Frazer 2006).

A related concern is how the effect of ratcheting is identified separately from the intertemporal substitution induced by the quota structure. The data are able to sort between these two separate dynamics in the following way. The extent of decline in the agent's observed sales after he crossed the ceiling in any quarter informs the model about the extent of intertemporal effort allocation induced by the quota structure. However, note that in the absence of ratcheting, effort, and hence, sales, should be strictly increasing between the quota and the ceiling. Hence, the extent of decline in the agent's observed sales after he crosses the quota, and before he attains the ceiling informs the model about the extent to which ratcheting plays a role. Figure 7 depicts the identification argument pictorially. The two other key parameters that are estimated in step 3 above are the cost (d) and risk aversion parameter (r). The cost of effort parameter is identified from the fact that sales are above the intercept

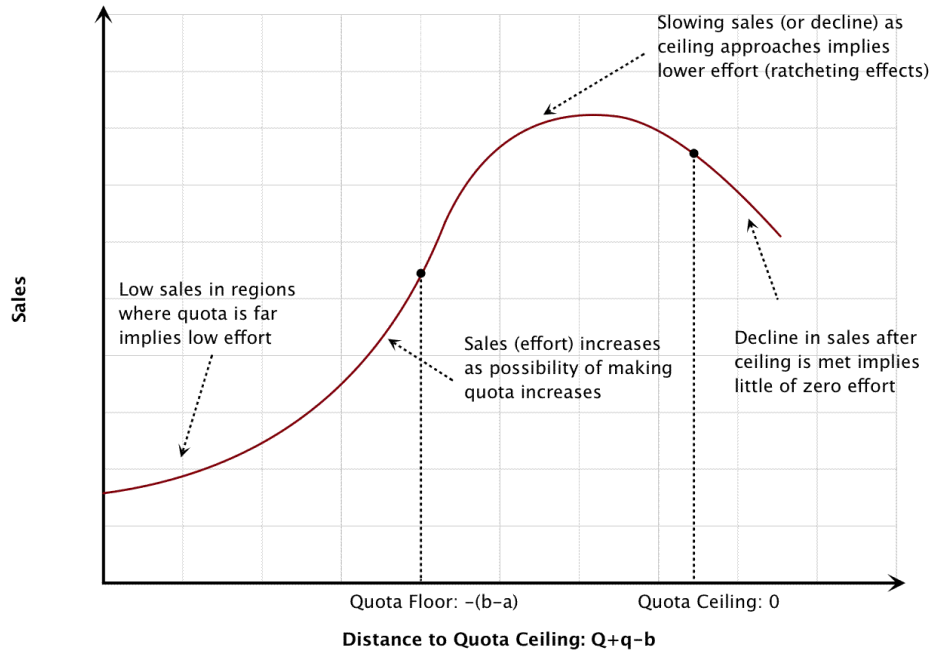


Figure 7: Identification of effort from sales profile.

in the first two months of the quarter. That is, if effort were costless, it would be optimal to exert no effort in the first two months and meet any target in the third month alone. The fact that effort is costly induces a constraint on how much sales can be generated in any given month. This, along with the structure of the sales response function, acts as the primary identification mechanism for the cost of effort parameter. Finally, the risk aversion parameter is identified by the degree to which effort (sales) changes due to changes in the variance of wealth. This variation in wealth is generated by within-agent factors that shift demand over time that are unrelated to the agent's distance to quota.

5 Data and Estimation Results

Table 8 presents summary statistics from our data. The sales-force has 87 salespeople who are about 43 years of age on average, and have been with the firm for approximately 9 years. The firm did not significantly hire, nor have significant employee

turnover in this sales-department during the time-period of the data.¹³ The average salesperson in the sales-force earns \$67,632 per annum via a fixed salary component. The annual salary ranges across the sales-force from around \$50,000 to about \$90,000. The firm’s output-based compensation is calibrated such that, on a net basis, it pays out a maximum of \$5,000 per agent per quarter, if the agent achieves 133% of their quarterly quota. On an average this implies that the agent has a 77%-23% split between fixed salary and incentive components if they achieve all targets. This is roughly what is achieved in the data. Across agents-quarters in the data, the average proportion of quarterly payout due to incentives is 16.8% (std. dev. 20.9%). Agents have exclusive territories and differ in terms of the number of clients they have, but are balanced in terms of the type of clients and the total sales-calls they are required to make.

The mean quota for the sales-force is about \$321,404 per quarter. The mean attained sales stands at \$381,210, suggesting that agents at the firm tend to target quarterly sales in the range in which incentives are earned. This is further evidenced by the fact that the range and dispersion parameters of the cumulative sales at the end of the quarter and the quota levels are also fairly close.

From Table (8), it appears on average that the firm adopts an asymmetric ratcheting approach to quota setting. When salespeople beat quotas the average increase in subsequent quarter quotas is about 10%, but on the flip side, falling short of quotas only reduces the next quarter quota by about 5.5%. This is consistent with some other earlier studies (e.g. Leone, Misra and Zimmerman 2004) that document such behavior at other firms, and is also consistent with our conversations with the firm management. Finally, the table documents that monthly sales average about \$138,149, a fairly significant sum.

5.1 Results from estimation

We now report the results from estimation. We first discuss the results from the first stage, which includes estimation of the effort policy function, and the quota transition process. Subsequently, we discuss the results from the estimation of the cost function and risk aversion parameters.

¹³So as to avoid concerns about learning-on-the job, and its interactions with quotas, 5 sales-agents, who had been with the firm for less than 2 years were dropped from the data.

5.1.1 Effort Policy

The effort policy function was estimated separately for each agent using a flexible Chebychev polynomial basis approximation. We approximate the effort policy using the tensor product of basis functions of dimension 2 in each of the two continuous state variables (cumulative sales and quota), allowing month specific intercepts, and allowing the first two basis functions to be month specific. We find that this specification fits the data very well. On average, we are able to explain about 79% of the variation in observed sales. Figure (8) plots a histogram of the R^2 values from the estimation across agents.

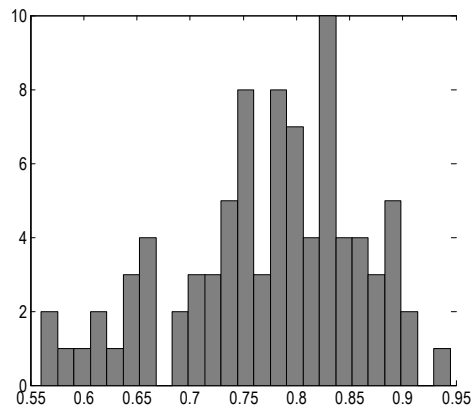


Figure 8: Histogram of R^2 values from first-stage effort policy estimation.

Rather than present estimates of the parameters of the basis functions approximating the effort policy, we present the estimates in graphical form. Figure 9 presents a first look at the effort policy using data pooled across all agents. The light areas in figure Figure 9 represent peaks of effort (dark areas representing valleys). Looking at Figure 9, we see that the data shows a clear pattern whereby effort tends to increase in the quota, which supports the “effort inducement” motivation for quotas noted by the theory. The variation of effort with cumulative sales is also intuitive. When cumulative sales are less than quota (areas to the left of the diagonal), the agent tends to increase effort. When cumulative sales are much greater than quota (areas to the right of the diagonal line), there is little incentive for the agent to exert further effort, and sales decline.

We now present contours of the effort policy estimated at the agent level. Figure

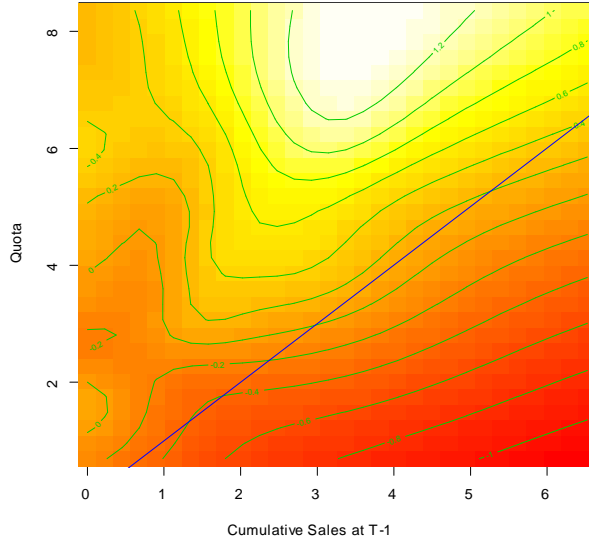


Figure 9: Contours of the Estimated Effort Policy.

10 shows the contours for nine of the sales-people. We find that there is considerable heterogeneity across the salespeople, which is consistent with wide variation in agent productivity. At the same time, we find that the basic pattern described above remain true. Similar to the average contour plot discussed below, we see sales increase with quota but fall after cumulative sales have exceeded quota.

5.1.2 Ratcheting Policy

We now discuss the estimated transition process for ratcheting. Figure (11) presents results from regressions in which we project the quota in quarter τ on flexible functions of agent's sales and quotas in quarter $(\tau - 1)$. Due to the fact that quotas vary only at the quarter-level, we estimate a pooled specification with agent fixed-effects. We are able to explain about 78% of the variation in quotas over time. Figure (11) also reports Breusch-Godfrey statistics for tests of 1st and 2nd order serial correlation in the ratcheting errors. Lack of serial correlation will imply that our flexible specification has captured all sources of persistence in the manager's quota updating policy. We see that with sufficient terms in the polynomial approximation, we are able to reject serial correlation in the ratcheting residuals.

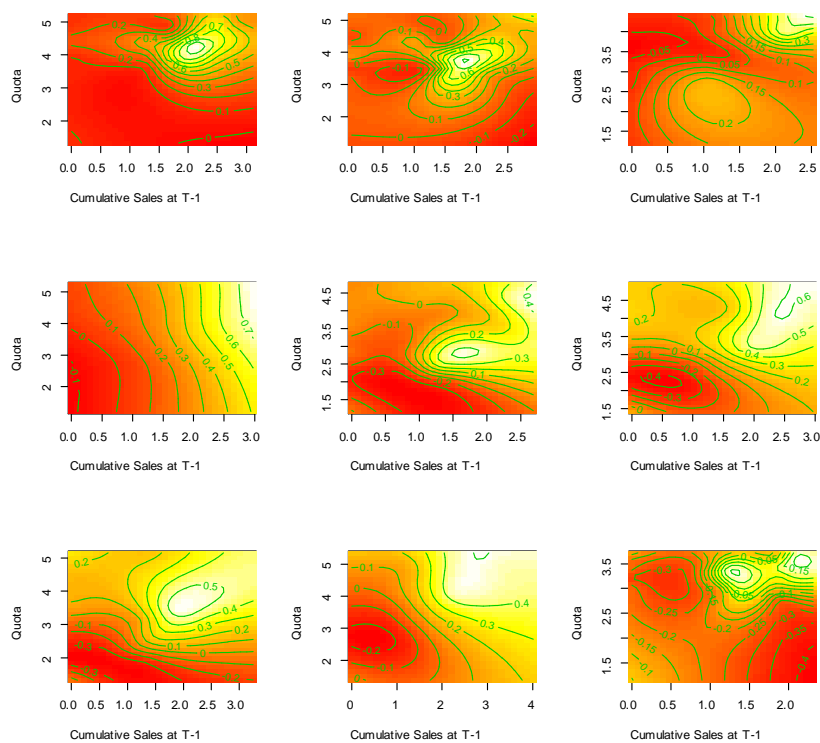


Figure 10: Examples of estimated effort policy functions across salespeople.

	Variable	Param	t-stat	Param	t-stat	Param	t-stat	Param*	t-stat*
	Constant	1.18	12.48	2.74	11.43	2.05	4.73	-0.02	-0.03
	$a(t-1)$	0.42	8.42	0.24	4.52	1.29	5.83	0.69	1.40
	$Q(t-1)$	0.32	6.09	0.17	2.92	-0.59	-2.57	1.42	7.76
	$a(t-1)^2$					-0.12	-4.89	-0.08	-0.75
	$Q(t-1)^2$					0.09	3.22	-0.33	-13.62
	$a(t-1)^3$							0.00	0.41
	$Q(t-1)^3$							0.04	30.21
	Agent fixed effects included?	N		Y		Y		Y	
	R^2	0.481		0.565		0.576		0.785	
	Breusch-Godfrey (1) p-value [#]	0.000		0.000		0.003		0.136	
	Breusch-Godfrey (2) p-value ^{##}	0.000		0.000		0.006		0.236	

*Preferred specification. Nobs = 1,044. Quotas and sales have been normalized to 100,000-s of \$\$s. [#]Tests against the null of zero 1st order serial correlation in the presence of a lagged dependent variable. ^{##}Tests against the null of zero 2nd order serial correlation in the presence of a lagged dependent variable.

Figure 11: Estimates of the Ratcheting Policy.

5.1.3 Second Stage Parameter Estimates

The remaining elements needed for the evaluation of counterfactual plans is an estimate for the joint distribution of the cost of effort (d) and risk aversion parameters (r). In this section we present estimates conditioned on the point estimates of μ_i . Figure (12) presents the estimated joint PDF of the two parameters.

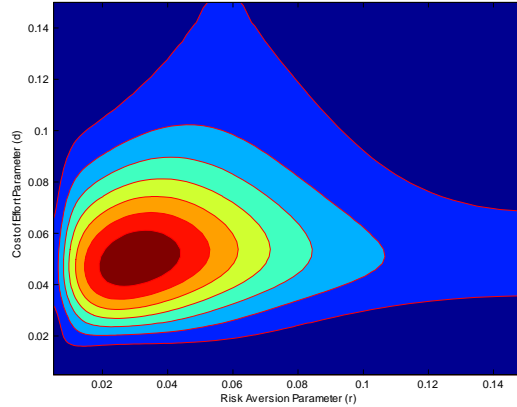


Figure 12: Joint PDF of the Cost of Effort and Risk Aversion parameters.

We find that there is a large amount of heterogeneity on both parameters. This is especially the case for the risk aversion parameter that varies significantly across agents from about 0.0018 to approximately 0.33. The large values are driven by the fact that for some salespeople, the variance of earnings across the sample period is low, resulting in risk aversion parameters that are correspondingly large. A more appropriate construct to examine is the monthly average risk premium, $\frac{r}{T} \mathbb{V}(W)$, which has a mean of around \$341.22 (median = \$281.36).

The density of the cost of effort parameter is much tighter with a mean of 0.0508 (median = 0.0471). There is still substantial heterogeneity in d as well with values ranging from 0.02 through 0.16. The parameters values translate approximately to a mean (across sales-agents) cost of effort of about \$1,591.7 per month. While not reported here, standard errors were computed using a bootstrap approach and are available from the authors upon request. The cost of effort parameter (d) was significant for all all agents at the 0.05 level while the risk aversion parameter (r) was significantly different from zero for 81 out of the 87 salespeople at the 0.05 level.

We also correlated the estimates of cost and risk aversion with the observed demographics of the sales-agents. Though we found that age and female gender correlated positively with risk, and tenure correlated negatively with the cost of effort, these were not statistically significant. To us, this underscores the importance of methods that flexibly control for unobserved heterogeneity in such contexts.

The characterization of $\mathcal{F}(d, r|\mu)$ completes the discussion of the estimation of various parameters needed for policy experiments. In what follows, we now discuss how we use these estimates, combined with the solution of the dynamic programming problem faced by the salesperson, to fine-tune the current compensation plan used by the firm. We then discuss results from the implementation at the firm of an improved plan based on the model.

6 Comparative Dynamics and a New Plan

In what follows we compare the current compensation policy at the focal firm to a series of alternative plans with the aim of uncovering causes and remedies of any distortions that may exist, which may then suggest more profitable configurations. The model provides us with estimates of the primitives underlying each sales agent's behavior. In principle, these primitives should allow us to search for the *optimal* compensation plan for the firm conditional on the estimated agent profile. Unfortunately, to our knowledge, there exists no straightforward algorithm that would implement an exhaustive search over the multidimensional compensation space and uncover the optimal second best compensation policy. As an alternate strategy, we conduct a number of counterfactual tests oriented toward evaluating the marginal profitability of the various compensation elements. While these comparative dynamics are conditioned on the particular characteristics of the current compensation plan, they allow us to investigate changes in agent behavior and output under the alternative compensation scenarios.¹⁴

We evaluate the different scenarios based on firms' expected profits under the new payment policy, as defined in equation (10).¹⁵ The empirical analog of (10) is

¹⁴As a caveat, note this is true as long as the alternate compensation schemes are not structurally different from the current plan. For example, relative compensation schemes, which condition compensation of a given agent on the performance of others, would require consideration of new elements such as fairness and competition which are not present in the current structure.

¹⁵To guard against the influence of outliers we integrate the profit function only over the interquartile range and renormalize the results.

constructed as follows:

$$\widehat{\Pi}_\varphi = \frac{1}{T \times NS} \sum_{s=1}^{NS} \sum_{\tau=0}^T \beta^\tau [q(e_\varphi; \Psi^s) - \wp(q(e_\varphi); \Psi^s)] \quad (29)$$

where $\wp(q(e); \Psi^s)$ is the compensation policy evaluated at a given draw Ψ^s from $\mathcal{G}_\varepsilon(\varepsilon_\tau) \times \mathcal{F}(\mu, r, d)$ and (e_φ) is the effort policy expended by the agent when faced with compensation policy $\wp(q(e_\varphi); \Psi^s)$,

$$e_\varphi = \arg \max_{e>0} V(s; e | \Psi^s, \wp(\cdot))$$

For our simulations we use $T = 25$ and $NS = 500$. In what follows below, all results at the monthly level refer to averages of expected profits or revenues over T .

We start by evaluating the three key features of the plan at the firm, viz., the ceiling, the quota, and the quota horizon. In addition, we evaluate the extent which better accommodation of heterogeneity in productivity across agents improves profits. We discuss the logic behind the changes we evaluated below.

(i) Removal of quotas and ceiling

As discussed earlier, the presence of the quota ceiling provides incentives to the salesperson to shade effort early in the quarter. One dynamic arises from the fact that the agent may find it optimal to wait to resolve uncertainty over the realization of early demand shocks, and to then make effort allocation decisions later in the quarter. This sequentially optimal strategy allows the agent to maximize the possibility of “making quota.” An additional dynamic arises from the expectation of quota ratcheting which may exacerbate the distortion by forcing salespeople to keep the realized output below the level of the ceiling. Both these effects are annulled by the removal of quotas. Hence, one counterfactual involves considering changes to the extent and the incidence of quotas in the plan.

(ii) Monthly compensation horizon

Another source of effort dynamics arises on account of the fact that the compensation horizon spans an entire quarter. If quotas are not high enough, agents may find it optimal to shirk in the early months, and make up sales later in the quarter. This may especially be relevant for the most productive agents. This suggest changing the length of quota horizon to improve the plan. Since agents

update their information sets at the end of each month (i.e. the institutional feature is that they access sales data only at the ends of each month), moving to a monthly plan would eliminate the effort shading problem. In other words, in a monthly plan, the agent can no longer wait for demand shocks to realize, and to then allocate effort, since the compensation period will have closed by then.

(iii) Heterogenous plans

Finally, our estimates suggest significant heterogeneity across agents. In particular, better fine-tuning of both salaries and incentives (commissions, and where applicable, quotas) based on this heterogeneity may increase profitability, since each plan would provide the right incentives to the agent based on their particular effort disutility and risk aversion profile.

We use the model to simulate profits as in equation (29) for each of the above scenarios. We find that changes in each of the described dimensions holding other aspects fixed, would improve profitability in each of the cases listed above. For the sake of brevity we do not outline each plan here but simply point out that the range of incremental profits ranges from 0.8% to 7.7% for these plans while revenues were predicted to increase between 2.3% and 13.4%. A complication arises on account of the complexity of the behavior induced by these changes. In particular, a finding of increased profits *ceteris paribus* does not imply that profits would increase if those changes are implemented jointly. While evaluating each permutation separately is impossible, our simulations generally reveal that joint consideration of the set of changes described above, is almost always found to increase profits. These plans were then proposed to the firm for their evaluation.

6.1 A New Plan: Implementation and Results

From a normative perspective, we believe there is value in discussing here the process by which options that may be considered superior by researchers on theoretical grounds, may be modified based on practical realities at the firm. A first-order issue arises because several cultural, legal and infrastructure constraints at the firm need to be accommodated prior to implementation; these constraints modify the set of plans that can be considered. For example, the firm in our data was not open to the idea of heterogeneous plans on account of a concern that it may engender a sense

of inequity and discontent amongst salespeople. Further, simple plans were valued more, as salespeople were clear about their dislike for complexity in compensation.¹⁶ These constraints narrowed the range of plans possible to a feasible set. The feasible set was then made available to the firm for consultation with various constituencies within the organization, including senior management, sales managers, salespeople and legal and human resources teams. A final plan was then chosen for implementation. This plan featured no quotas or ceilings, and a monthly incentive based on a straight commission on sales. Due to confidentiality concerns, we cannot divulge the exact value of the commission nor further details of its implementation.

The simulations from our model predict that the firm's revenues would increase by about 8.2%, and profits, measured as revenues minus compensation payouts, would increase by about 5.1% under this new plan. Our simulations suggest that the impact of the new plan will be quite varied across the salesforce. Figure 13 shows the predicted impact of the new plan across the salesforce. While the majority of the salesforce exhibits improvements in the 0-10% range there are some salespeople for whom the change is expected to be quite large and others for whom there is even a drop in sales. We should point out here that this plots represents the average across a number of simulations and that the actual impact of the new plan could deviate from this particular pattern pattern. The reader should note that these estimates do not reflect the fixed costs of implementation such as changes to the HR. information technology system, sales-force retraining and other transition costs. Implicitly, we assume these are the same under all considered options.

Once the final scheme was chosen, the firm implemented a transition plan which included educating salespeople about the changes to their compensation and the managing their expectations. For example, detailed worksheets were provided which showed the compensation the agents would receive under the new plan at various points of performance. Particular care was taken to assuage fears and risk perceptions related to the change and salespeople were engaged in numerous discussions to ensure their doubts and questions were answered.¹⁷ The new plan was implemented January 1st 2009. In the remainder of this section, we present results on the sales performance under the new plan using data from the first two quarters of the 2009 calendar year.

¹⁶We consider the valuation by agents of simplicity, and its manifestation as menu costs to the firm, an important direction for future research.

¹⁷To assuage concerns about this period of change in the firm, the data from the last two quarters of 2008 are not used in the estimation of the model.

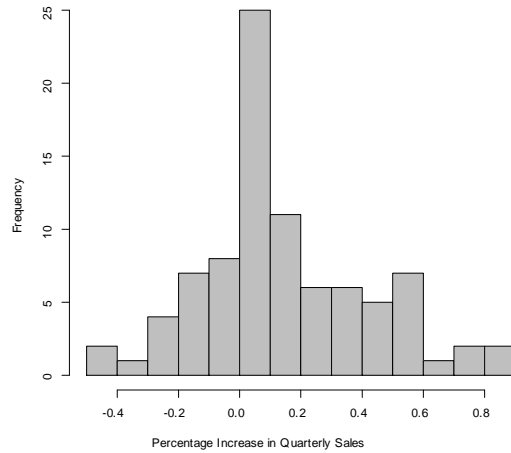


Figure 13: Predicted percentage change in quarterly revenues under new plan relative to current plan.

6.1.1 Results

Aggregate Effects We start by discussing aggregate impact of the new compensation scheme. The aggregate effect is about a 9% increase in revenues across agents. Figure 14 shows the distribution of the percentage and dollar changes in revenues in the first quarter of 2009 relative to the first quarter of 2008 at the agent-level. Looking at Figure 14, we see that in percentage terms, the new plan provided a lift in revenues at the agent-level of about 22.6% (std. 18.8%) on average. In dollar terms, this translates to about \$79,730 per agent per quarter on average (std. \$62,809). Importantly, the fact that overall quarterly sales increased suggests that the old plan was *inefficient*. In particular, the fact that output went up indicates that dynamics under the old plan did not simply have the effect of shifting the timing of doctor’s prescriptions, but rather, also reduced the aggregate long-run orders from doctors. This is consistent with a story where doctors simply prescribe substitute contact lenses from other brands when agents respond with low effort arising from incentives. Later in this section, we provide further evidence that the old plan accentuated such brand-switching by doctors, and that the new plan reduces this inefficiency.

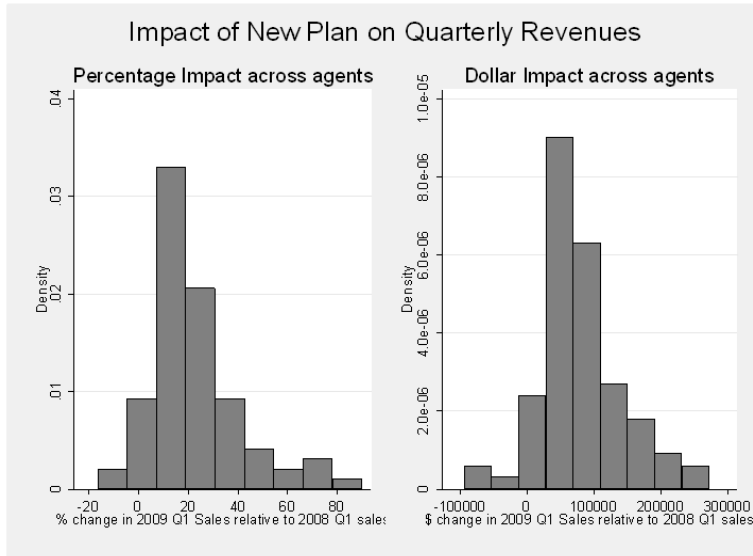


Figure 14: Quarter level effects of the new plan.

Distributional Effects From Figure 14, we also see that there is considerable heterogeneity in the extent to which agents respond to the new plan. While some agents reduce effort, most agents seem to have increased effort and output. The pattern of the heterogeneity is not very different from that obtained from our model and presented in Figure 13.

We now discuss whether the estimates from our structural model are indicative of the extent to which agents may respond to the new plan. We use the nature of response of agents to the new plan to informally assess the face validity of the estimates we obtained under our maintained assumptions of agent behavior. In particular, theory suggests that agents with higher cost of effort will find it harder to achieve sales under the new plan. Hence we would see sales under the *new plan* to be lowest for agents with high cost of effort. To assess this we ran a simple regression of sales under the new plan as a function of month of quarter and the cost of effort and risk aversion parameters. The results of these regressions were striking. First, the regression as a whole was significant (p -value of 0.016 for the F -statistic). Second, we find that *both* the risk aversion coefficient and the cost of effort parameter had negative and significant effects on sales (p -values = [0.0414, 0.0002]). These results point to the fact that the estimated structural parameters are indeed able to correlate to the observed behavior in the field. Lastly, the regression also indicated that

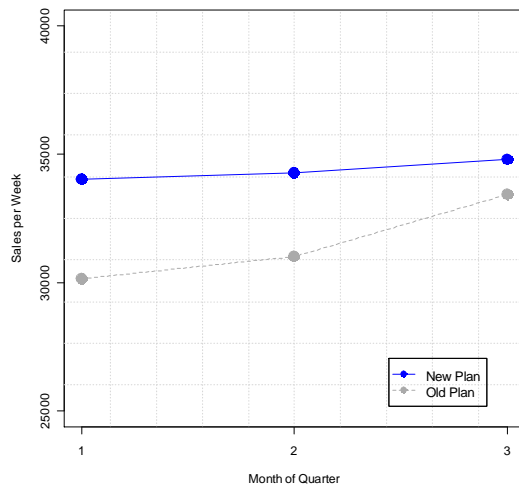


Figure 15: Sales patterns under the old and new plans.

the month of quarter had no significant impact on sales. This is important and we discuss this in more detail in what follows.

Month-level Effects We now use the data from the new plan to assess whether the within-quarter changes in monthly sales follow patterns we predicted based on our model. In particular, we use the data under the new plan to assess the importance of the two main alternative explanations for the intertemporal patterns in sales observed previously, viz. stockpiling and seasonality. Under the new plan, incentives are the same across months. Hence, under the null of no stockpiling or seasonality, we should see that sales are flat across months of the quarter. Figure (15) shows a plot of sales-per-week across months of the quarter for both the old and the new plan. The plot from the old plan replicates the “scallop” pattern suggestive of the inefficiency. The sales under the new plan is found to be flat, corroborating the model-free analysis earlier that ruled out seasonality and stockpiling.

We also investigate the shifts in sales that occurred under the new plan. Figure (16) plots the kernel density of the percentage change in sales across agents across months of the quarter for Q1-2009 relative to Q1-2008. Following Figure (16), we see that relative to the old plan, month 3 sales have shifted down, and month 1 sales have shifted upwards, which is consistent with the finding from the previous empirical analysis that there is likely shirking in the early months of the quarter under the old

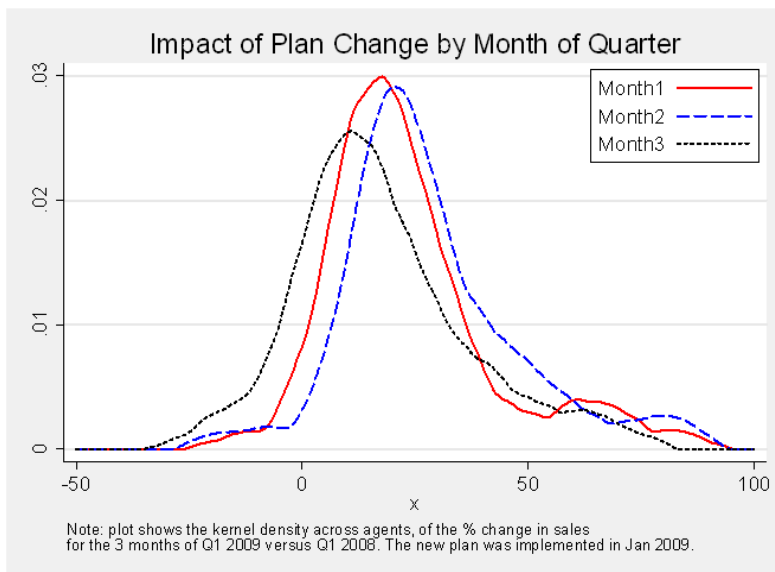


Figure 16: Month level effects of the new plan.

plan.

Another look, and client-level effects We close our discussion of the new plan by exploring changes in output at the client level. We start by reporting on regressions of sales-per-week at the agent-month level. These results are presented in Figure (17). We first look only at the first two-quarters of 2009 (column [1]), and run a regression of sales per-week, per-agent on month-fixed effects, and test whether the month-fixed effects are significantly different from zero. Essentially, we ask whether sales are flat within the quarter under the new plan. Figure (17) reports on the p -value from this test: the null that the sales-per-week is the same across months of the quarter is not rejected. This corroborates the rejection of seasonality/stockpiling from the previous section. For contrast, we report the analogous regression using all the data from the previous plan. Looking at column [2], we see that the null that month effects are the same is strongly rejected (p -value = 0.000).

The first regression in Column [3] now reports on the effect of the new plan by pooling data across pre- and post-plan implementation months. From column [3], we see that the incremental effect of the new plan is to add about \$2,827 per agent per week on average. This translates to an incremental +\$11,308 ($\2827×4) for each

agent-month. If we multiply by 87 agents, we get a figure of about \$0.983 million incremental per month company-wide, which is about \$12 million annually. While profits are trickier to nail down, numbers we have obtained from the firm suggest that overall profitability after the implementation of the plan has increased by over 6%. The second set of regressions in column [3], now splits the main effect of the new plan by month-of-the-quarter. Consistent with the previous analysis, we see that relative the old plan, sales in month 3 has reduced, while sales in month 2 remains the same.

Finally, to translate these effects to the client-level, from Table (8), we note that there are 162 clients per agent. Hence, the incremental effect of the new plan translates to \$69.81 ($\$11,308/162$) per doctor per month. This is roughly 1-2 prescriptions (average cost of contact lens is about \$35 to \$50 per box), which is small at the individual physician-level. That the effect at the individual-physician level is small enough also suggests that competitive reaction to the change in the plan is likely muted; hence, it seems reasonable to interpret the numbers reported here as the long-run effect of the improvements in the compensation scheme.

Variable	Data New Plan Only [1]		Data for Old Plan Only [2]		All Data [3]	
	Param	t-stat	Param	t-stat	Param	t-stat
Constant	\$ 34,015.2	68.56	\$ 30,161.0	93.26	\$ 30,389.2	105.59
Month 2	\$ 250.4	0.36	\$ 861.5	1.88	\$ 725.7	1.87
Month 3	\$ 789.5	1.13	\$ 3,259.6	7.13	\$ 2,710.7	6.97
New Plan					\$ 2,827.2	7.40
(New Plan)*(Month 2)					\$	(611.1)
(New Plan)*(Month 3)					\$	(2,470.2)
pvalue(H0: Month 2 = Month 3 = 0)	0.5167		0.0000			
R2	0.0027		0.0308		0.0461	
Nobs	492		1722		2214	
					0.0494	2214

Notes: dependent variable is sales-per-week at the agent-month level. The new plan was implemented in Jan 2009, and sales for each agent under the new plan are available for the 1st six months of 2009. [1] uses data for the first six months of 2009; [2] uses data for the 1st six months of 2005-2008. [3] uses data for the 1st six months of 2005-2009, 2009 inclusive.

Figure 17: Client-level effects of the new plan.

Epilogue The firm also reports that employee satisfaction with the new plan is high, arising primarily from the elimination of quotas, and the associated subjective evaluation induced via ratcheting. Overall, our results suggest that the new plan has been a success. Further, the results support the external validity of the model and the estimates, and strongly support the validity of dynamic-programming based agency-theory models for assessing and improving real-world sales-force compensation schemes.

7 Conclusions

This paper presented a comprehensive framework to analyze and fine-tune sales-force compensation schemes. The framework is built on agency theory, and is implemented using numerical dynamic programming. The framework is flexible enough to handle nonlinearities and kinks commonly observed in real-world contracts. The framework emphasizes the careful consideration of the dynamics induced in agent behavior by these aspects of compensation schemes. An algorithm for estimating the parameters indexing the model is also proposed. The algorithm places a premium on flexible, nonparametric accommodation of unobserved heterogeneity, and exploits the richness of informative, internal firm-databases linking contracts and output. The external validity of the framework is demonstrated via a field-implementation at the company that provided the data. The field implementation increases revenues substantially. Further, patterns of changes in sales are found to be consistent with the predictions from the model, and validates the assumptions employed.

We wish to conclude by discussing caveats and possible extensions. An important caveat is that the framework is not intended to be applied to durable goods which exhibit buyer-side demand dynamics, or to goods with buyer-side seasonality. More data that enables controls for these aspects would be needed in those contexts. The framework will have to be extended to consider plans that are structurally different from the one addressed here (e.g. relative performance schemes or tournaments). Plans that result in dependencies across agents are especially complex, and require an additional equilibrium concept for solution.

Several extensions are possible. Computing the optimal plan is an unsolved, but methodologically challenging endeavor. Accommodating potential multi-tasking by agents is another important area for future research. Finally, better estimation

of agent primitives, especially agent's discount factors, will help better pin down the key dynamics we describe. New methods proposed recently for measuring discounting (e.g. Dube, Hitsch and Jindal 2009), thus hold great promise for analyzing sales-force compensation.

8 References

1. Albers, S. and Murali Mantrala (2008), "Models for Sales Management Decisions," Handbook of Marketing Decision Models.
2. Akerberg D., K. Caves and G. Frazer (2006). "Structural Identification of Production Functions," working paper, UCLA.
3. Arcidiacono, P. and Bob Miller (2008), "CCP Estimation of Dynamic Discrete Choice Models with Unobserved Heterogeneity," working paper, Duke University.
4. Asch, B. (1990). "Do Incentives Matter? The Case of Navy Recruiters," Industrial and Labor Relations Review, 43 (3): 89S-106S.
5. Bajari, P., C. L. Benkard and J. Levin (2007). "Estimating Dynamic Models of Imperfect Competition," Econometrica, 75(5), 1331-1370.
6. Baker, G., R. Gibbons and K. Murphy (1994). "Subjective Performance Measures in Optimal Incentive Contracts," Quarterly Journal of Economics, 109 (4): 1125-1156.
7. Basu, A., R. Lal, V. Srinivasan and R. Staelin (1985). "Sales-force Compensation Plans: An Agency Theoretic Perspective," Marketing Science, 8 (3): 324-342.
8. Bhardwaj, P. (2001). "Delegating Pricing Decisions," Marketing Science, 20 (2): 143-169.
9. Chevalier, J. and G. Ellison (1997). "Risk Taking by Mutual Funds as a Response to Incentives," Journal of Political Economy, 105 (6): 1167-2000.
10. Cho, S. and J. Rust (2008), "Is Econometrics Useful for Private Policy Making? A Case Study of Replacement Policy at an Auto Rental Company," Journal of Econometrics 145 243-257.
11. Chung, D., Thomas Steenburgh and K. Sudhir (2009), "Do Bonuses Enhance Sales Productivity? A Dynamic Structural Analysis of Bonus-Based Compensation Plans," working paper, Yale University.
12. Copeland, A. and Monnett, C. (2009), "The Welfare Effects of Incentive Schemes," Review of Economic Studies, 76, pp. 96-113.

13. Coughlan, A. and C. Narasimhan (1992). "An Empirical Analysis of Sales-force Compensation Plans," *Journal of Business*, 65 (1): 93-121.
14. Coughlan, A. (1993). "Sales-force Compensation: A Review of MS/OR Advances," *Handbooks in Operations Research and Management Science: Marketing* (vol. 5), Gary L. Lilien and Jehoshua Eliashberg, editors, Amsterdam: North-Holland.
15. Coughlan, A. and S. Sen (1989). "Sales-force Compensation: Theory and Managerial Implications," *Marketing Science*, 8(4), 324-342.
16. Courty, Pascal and Jerry Marschke, (1997), "Measuring Government Performance: Lessons from a Federal Job-Training Program," *American Economic Review*, 87, pp. 383-88.
17. Dube, J-P, Hitsch, G. and P. Jindal (2009), "Estimating Durable Goods Adoption Decisions from Stated Preference Data," working paper, University of Chicago.
18. Godes, D. (2003). "In the Eye of the Beholder: An Analysis of the Relative Value of a Top Sales Rep Across Firms and Products," *Marketing Science*, 22 (2): 161-187.
19. Healy, P. (1985). "The Effect of Bonus Schemes on Accounting Decisions," *Journal of Accounting and Economics*, (7) 1-3: 85-107.
20. Holmstrom, B. (1979). "Moral Hazard and Observability," *Bell Journal of Economics*, 10: 74-91.
21. Holmstrom, B. and P. Milgrom (1987). "Aggregation and Linearity in the Provision of Intertemporal Incentives," *Econometrica*, 55, 303-328.
22. Jiang, Renna and R. Palmatier (2009). "Structural Estimation of a Moral Hazard Model: An Application to Business Selling," working paper, U.C. Davis School of Management.
23. Joseph, K. and M. Kalwani (1992). "Do Bonus Payments Help Enhance sales-force Retention?" *Marketing Letters*, 3 (4): 331-341.
24. Lal, R. and V. Srinivasan (1993). "Compensation Plans for Single- and Multi-Product sales-forces: An Application of the Holmstrom-Milgrom Model," *Management Science*, 39 (7):777-793.
25. Larkin, I. (2006). "The Cost of High-Powered Incentive Systems: Gaming Behavior in Enterprise Software Sales," working paper, Harvard Business School.
26. Lazear, E. (1986). "Salaries and Piece Rates," *Journal of Business*, 59 (3): 405-431.

27. Lazear, E. (2000). "Performance, Pay, and Productivity," *American Economic Review*, 90 (5), 1346-1361.
28. Lee, Donald and Zenios, Stefanos (2007), "Evidence-Based Incentive Systems With an Application in Health Care Delivery," working paper, Stanford Graduate School of Business.
29. Leone A., S. Misra and J. Zimmerman (2004). "Investigating Quota Dynamics", working paper, University of Rochester.
30. Li, T. and Q. Vuong (1998), "Nonparametric Estimation of the Measurement Error Model Using Multiple Indicators," *Journal of Multivariate Analysis*, Vol. 65, No. 2, 139-165.
31. Mantrala, M., P. Sinha and A. Zoltners, (1994). "Structuring a Multiproduct Sales Quota-Bonus Plan for a Heterogeneous Sales Force: A Practical Model-Based Approach" *Marketing Science*, 13(2), 121-144.
32. Mantrala, M., P.B. Seetharaman, Rajeev Kaul, Srinath Gopalakrishna & Antonie Stam (2006), "Optimal Pricing Strategies for an Automotive Aftermarket Retailer," *Journal of Marketing Research*, 43, 4, 588-604.
33. Misra S., A. Coughlan and C. Narasimhan (2005). "Sales-force Compensation: An Analytical and Empirical Examination of the Agency Theoretic Approach," *Quantitative Marketing and Economics*, 3(1), 5-39.
34. Nair, H. (2007). "Intertemporal Price Discrimination with Forward-looking Consumers: Application to the US Market for Console Video-Games," *Quantitative Marketing and Economics*, 5(3), 239-292.
35. Olley S. and A. Pakes (1996). "The Dynamics of Productivity in the Telecommunications Equipment Industry," *Econometrica*, 64(6), 1263-1297.
36. Oyer, P. (1998). "Fiscal Year Ends and Nonlinear Incentive Contracts: The Effect on Business Seasonality," *Quarterly Journal of Economics*, 113 (1): 149-185.
37. Oyer, P. (2000). "A Theory of Sales Quotas with Limited Liability and Rent Sharing," *Journal of Labor Economics*, 18 (3), 405-426.
38. Prendergast, C. (1999). "The Provision of Incentives in Firms," *Journal of Economic Literature*, 37 (1): 7-63.
39. Raju, J. S., and V. Srinivasan (1996). "Quota-based compensation plans for multi-territory heterogeneous sales-forces," *Management Science* 42, 1454-1462.
40. Rao, R., (1990). "Compensating Heterogeneous Sales-forces: Some Explicit Solutions," *Marketing Science*, 9(4), 319-342

41. Rust, J. (1996), "Numerical Dynamic Programming in Economics," Handbook of Computational Economics, Chapter 14.
42. Steenburgh, T. (2008). "Effort or Timing: The Effect of Lump-sum Bonuses," Quantitative Marketing and Economics, 6:235-256.
43. Weitzman, M. (1980), "The 'Ratchet Principle' and Performance Incentives," Bell Journal of Economics 11, 302-308.
44. Zenger, T. and S. Lazzarini (2004). "Compensation for Innovation: Do Small Firms Offer High-Powered Incentives That Lure Talent and Motivate Effort?," Managerial and Decision Economics, 25: 329-345.
45. Zoltners, A., P. Sinha and G. Zoltners (2001). "The Complete Guide to Accelerating Sales Force Performance," American Management Association, New York.
46. Zoltners, A., Prabhakant Sinha and Sally E. Lorimer (2008). "Sales Force Effectiveness: A Framework for Researchers and Practitioners," Journal of Personal Selling and Sales Management, 28 (2), 115-131.

Table 1: Descriptive Statistics of Data.

Variable	Mean	SD	Min	Max
Agent Demographics				
Salary	\$67,632.28	\$8,585.13	\$51,001.14	\$88,149.78
Incentive Proportion at Ceiling	0.23	0.02	0.8	0.28
Age	43.23	10.03	27	64
Tenure	9.08	8.42	2	29
Number of Clients	162.20	19.09	63	314
Quarter Level Variables (across agents)				
Quota	\$321,404	\$86,112.67	\$178,108.93	\$721,770.14
Cumulative Sales (end of quarter)	\$381,210	\$89,947.66	\$171,009.11	\$767,040.98
Percent Change in Quota (when positive)	10.01%	12.48%	00.00%	92.51%
Percent Change in Quota (when negative)	-5.53%	10.15%	-53.81%	-00.00%
Monthly Level Variables (across agent-months)				
Monthly Sales	\$138,149	\$383,19.34	\$45,581.85	\$390,109.07
Cumulative Sales (beginning of month)	\$114,344	\$985,94.65	\$0	\$65,2474.25
Distance to Quota (beginning of month)	\$278,858	\$121,594.2	\$20,245.52	\$83,5361.10
Number of Salespeople	87			

A Appendix A: Computational Details

This appendix provides computational details of solving for the optimal policy function in equation (9) and for implementing the BBL estimator in equation (28).

Solution of Optimal Policy Function The optimal effort policy was solved using modified policy iteration (see, for e.g., Rust 1996 for a discussion of the algorithm). The policy was approximated over the two continuous states using 10 points in each state dimension, and separately computed for each of the discrete states. The expectation over the distribution of the demand shocks ε_t and the ratcheting shocks v_{t+1} were implemented using Monte Carlo integration using 1000 draws from the empirical distribution of these variates for the agent. The maximization involved in computing the optimal policy was implemented using the highly efficient SNOPT solver, using a policy tolerance of 1E-5.

Estimation of Agent Parameters We discuss numerical details of implementing the BBL estimator in equation (28). The estimation was implemented separately for each of the 87 agents. The main details relate to the sampling of the initial states, the generation of alternative feasible policies, and details related to forward simulation. For each, we sampled a set of 1002 initial state points uniformly between the minimum and maximum quota and cumulative sales observed for each agent, and across months of the quarter. At each of the sampled state points, we generated 500 alternative feasible policies by adding a normal variate with standard deviation of 0.35 to the estimated optimal effort policy from the first stage (effort is measured in 100,000-s of dollars). Alternative feasible policies generated by adding random variates with large variances (e.g. 5), or by adding noise terms to effort policies at only a small subset of state points, were found to be uninformative of the parameter vector. At each sampled state point, we simulated value functions for each of the 500 alternative feasible policies by forwards-simulating the model 36 periods ahead. The sample analog of the moment conditions are then obtained by averaging over the sampled states and alternative policies.