

Customer Acquisition, Engagement, and Retention in Online Advertising

A dissertation presented

by

Michael Els

to

The Harvard Business School

in partial fulfillment of the requirements

for the degree of

Doctor of Philosophy

in the subject of

Business Administration

Harvard University

Cambridge, Massachusetts

October 2019

©2019 Michael Els
All rights reserved.

Dissertation Advisor:
Professor Sunil Gupta

Author:
Michael Els

Customer Acquisition, Engagement, and Retention in Online Advertising

Abstract

Online advertising continues to evolve at a rapid as the internet and the digital marketing landscape mature. Firms face new challenges in acquiring, engaging and retaining customers. Entire new markets and technologies have grown out of the race digital marketing dominance. This dissertation aims to examine some of these advances and offer practical insights for today's firms that need to navigate this new world.

In the first essay, I explore the effects of user attention to online display advertising. Using two observational studies, I show that attention is highly heterogeneous and predictable during the user browsing session. The implications are that publishers should be more selective in ad placement and that advertisers should be more selective in ad purchases.

The second essay examines how programmatic advertising firms should efficiently allocate ads in real-time bidding environments on behalf of their client advertisers. I introduce the demand side platform problem which is related to both the adwords and publisher problems, but distinct in that the supply of ad space assumed unlearnable. I provide a real-time mechanism for efficient ad allocation in this setting and demonstrate efficacy using real-time bidding data.

In the final essay, I examine cross-merchant spillovers in coalition loyalty programs. I examine a natural experiment where a large grocery store joined a large loyalty program coalition. Using a

quasi-difference-in-difference approach and Bayesian Structural Time Series for causal inference, I find that adding a large complementary merchant into a coalition loyalty program increases sales and purchase frequency of existing customers at existing merchants.

Contents

Abstract.....	iii
Acknowledgements.....	vii
Introduction.....	1
1 Online Task Progression and Display Ad Engagement.....	3
1.1 Introduction.....	3
1.2 Literature Review.....	5
1.3 Attention to Tasks.....	7
1.4 Conceptual Model.....	9
1.5 Measurement.....	11
1.6 Design and Identification.....	13
1.7 Data Set 1 – Single News Tasks.....	18
1.7.1 Data Description.....	18
1.7.2 Empirical Model.....	20
1.7.3 Results.....	22
1.7.4 Propensity Score Matching Design.....	24
1.7.5 Results of Propensity Score Matching.....	25
1.8 Data Set 2 – Browsing Sessions with Multiple Tasks.....	31
1.8.1 Data description.....	32
1.8.2 Propensity Score Matching Design.....	34
1.8.3 Results.....	35
1.9 Validation Simulation for Advertisers.....	38
1.10 Conclusion.....	44
1.10.1 Caveats.....	45
1.10.2 Implications of Findings.....	45
2 Real-time Digital Ad Allocation: A Fair Streaming Allocation Mechanism.....	47
2.1 Introduction.....	47
2.2 Literature Review.....	49
2.3 Fairness.....	53
2.4 Algorithm Design.....	55
2.4.1 Environment.....	55

2.4.2	Proposed Solution	58
2.4.3	Part 1 - Agent level optimization	58
2.4.4	Market level allocation	65
2.5	Evaluation.....	69
2.5.1	Single-Advertiser Simulation.....	70
2.5.2	Multiple Advertisers (Market-level evaluation)	75
2.6	Conclusion.....	80
3	Cross-Merchant Spillovers in Coalition Loyalty Programs	82
3.1	Introduction	82
3.2	Related literature	85
3.3	Background, data, and preliminary evidence	87
3.4	Matching Analysis.....	92
3.5	Bayesian Structural Time Series	102
3.6	Extension to other merchants	112
3.7	Conclusion, Discussion, and Limitations.....	115
	References.....	118
	Appendix.....	124
A	Online Task Progression and Display Ad Engagement.....	124
A.1	Propensity Score Matching Details	124
A.2	SEC DSP Cost Estimates	125
B	Real-time Digital Ad Allocation: A Fair Streaming Allocation Mechanism	125
B.1	PID results without PSO	125
B.2	Market level DH formulation.....	127
C	Cross-Merchant Spillovers in Coalition Loyalty Programs	128
C.1	Variable Definitions	128
C.2	Matching Analysis.....	130
C.3	Other merchant entries	132
C.4	Market Level Bayesian Structural Time Series.....	136

Acknowledgements

I would like to thank my advisors, Sunil Gupta, Donald Ngwe, David Parkes, and Thales Teixeira for their guidance throughout this process. Their mentorship and instruction has been life changing and I will forever be indebted to them. I would also like to thank Ayelet Israeli, Elie Ofek, Rajiv Lal, and many other members of the Marketing unit who helped me tremendously.

I am incredibly grateful to my friends and family for their support and encouragement throughout my studies.

Introduction

Digital marketing is rapidly evolving and continues to upend the strategic marketing landscape. This is primarily due to the shift in consumer habits where people now spend a large proportion of their waking hours online. Unlike traditional media channels, digital channels provide a wealth of user data, instant feedback, and interactive contact. However, much of the digital space is still unsophisticated leading to cluttered websites and poor user engagement. Despite the abundance of online ads, richness in targetability data, and opportunity to engage users, firms have struggled to improve marketing efficiency in these new channels.

One area of rapid innovation is the movement of digital advertising to programmatic bidding. This has created a new ecosystem of firms that transact, track, and target users at the impression level. Advertisers now need to evaluate individual digital ads spaces in under a second and optimally place their ads in a wildly unpredictable environment.

Another area of innovation is experimentation in digital loyalty programs. These loyalty programs are also becoming an essential part of mobile advertising and modern sales promotions as marketers continue to find new ways to engage customers. These innovations continue to advance at a rapid pace despite the lack of evidence to support the formation of digital loyalty programs.

This dissertation adds to the digital advertising and CRM literature through empirical analysis of these emerging trends. In the first essay, co-authored with Thales Teixeira, I explore the effects of user attention to online display advertising. Much of previous research jumps ahead to later funnel activities assuming that all display ads are equally likely to be attended to or that attention is randomly distributed throughout time. I show that attention is highly heterogeneous

and predictable during the user browsing session. The implications for publishers are that many of the ad spaces on their websites are focused in areas where people will not attend to them. It also has implications for advertisers targeting practices in that they should not purchase ad space that is unlikely to ever be attended to.

The second essay, co-authored with Sunil Gupta and David Parkes, examines how programmatic advertising firms should efficiently allocate ads in real-time bidding environments on behalf of their client advertisers. I introduce this as the demand side platform problem which I compare to both the adwords and publisher problems. These firms face a highly unstable supply of ads that do not follow a stationary distribution and therefore violating the assumptions of the adwords and publisher problems. Additionally, these firms also face ad quota constraints from their client advertisers. Given the nature of the environment, existing parametric solutions are not feasible, and I provide a real-time mechanism for efficient ad allocation.

In the final essay, co-authored with Sunil Gupta and Donald Ngwe, I turn from acquiring customers to retaining them. I examine cross-merchant spillovers in coalition loyalty programs. This essay adds to the literature on loyalty program efficacy by examining a natural experiment where a large grocery store joined a large loyalty program coalition. I conduct our analysis using a quasi-difference-in-difference approach and Bayesian Structural Time Series for causal inference. I find that adding a large complementary merchant into a coalition loyalty program increases sales and purchase frequency of existing customers at existing merchants.

1 Online Task Progression and Display Ad Engagement

1.1 Introduction

Display advertising is ubiquitous on the Internet despite limited evidence of its efficacy. Not knowing where to place online ads that will garner attention, many marketers resort to overusing them. Marketing research has improved our understanding of when online ads work to some extent. However, much of the academic research in online advertising has focused on measuring user heterogeneity or modeling purchase behavior to better segment and target consumers (Manchanda et al. 2006, Rutz and Bucklin 2012, Hoban and Bucklin 2015). Past research prior to the advent of the Internet has shown the critical importance of timing in influencing consumers through advertising (Strong 1977). However, little has been done on understanding what consumers are doing online when exposed to internet advertising and how this impacts ad effectiveness. This paper incorporates the role of timing—as it pertains to the consumer’s progression through tasks routinely done online—on engagement to online display ads.

A critical factor in display advertising is the availability of attention at the time of exposure. We propose that consumer tasks and their progression through these tasks are fundamental to explain engagement to display advertisements. We build on current models of attention from psychology and marketing to develop a conceptual model of attentional buildup followed by attentional release during typical consumer browsing tasks. We hypothesize that consumers are more likely to attend to online display ads at the beginning and ending of tasks, as opposed to in the middle portion.

We test our hypothesis using two datasets. The first contains video ad view-through rates from experimental data randomized across users, ad campaigns and time. The second contains display ad click rates on a larger and richer observational dataset from a digital advertising targeting firm. Both datasets offer impression level measures of user engagement allowing for precise and high-powered identification. For each of these datasets, we estimate a logistic regression model with random effects to measure the impact of task progression while accounting for the standard predictors of display ad engagement. Furthermore, to overcome observational biases and endogeneity concerns in our data, we test our hypothesis using a novel application of Rubin's (1974) propensity score matching procedure for causal inference. Our implementation of this approach exploits website content and hierarchical page structure to predict task progression using the random forest algorithm. Our matching procedure utilizes binary search trees which enables our analysis to scale to big data.

Using these models, we find that an individual's task progression has a substantial effect on display ad engagement. We provide strong support for the hypothesis that consumers are more likely to attend to online display ads at the beginning and ending of tasks, as compared to in the middle. Our findings are conceptually and practically distinct from both the industry belief that ads at the top of web pages will have higher attention levels (Brebion 2018), as well as from prior research showing that ad engagement is equally likely throughout the browsing session (Chatterjee et al. 2003).

Our findings on ad engagement and task progression have important managerial implications for both web publishers and advertisers seeking to optimize online ad placements. One such implication is that, akin to how advertisers should understand how people search online for products in order to better determine keywords to bid for search engine ads, advertisers should

also understand people's primary tasks and their progression in order to better place display ads. Additionally, we conduct an out-of-sample simulation to demonstrate how advertisers can incorporate task progression in better targeting digital advertising.

This paper is organized as follows. We first conduct a brief literature review on display ad placement and attention, followed by formalizing our hypothesis of attention as a function of display ad location during tasks. We then develop our conceptual model and discuss our design and identification strategy to test our hypothesis. This is followed by the empirical estimation using two data sets showing support for our hypothesis and a simulation showing the practical benefits. In the final section, we conclude and offer recommendations for the placement of online display advertising.

1.2 Literature Review

The importance of digital media has led to significant research into digital advertising. Most early research on digital advertising has focused on understanding the effects of display advertising at the awareness stage by measuring changes in brand awareness, brand attitudes and purchase intentions (Ilfeld and Winer 2002, Dreze and Hussherr 2003, Cho and Cheon 2004, Moore et al. 2005). More recent research has instead focused on actual purchase outcomes from display advertising and factors such as privacy or segmentation strategies that affect later funnel purchase behavior (HobanBucklin 2015). However, little research has been done examining consumer attention to online tasks and how attention to tasks affects engagement with display ads. If attention is the bottleneck to downstream effects as those early studies showed, then the impact of browsing progression on reducing the attention available for display ads is a possible culprit for lack of ad effectiveness.

One paper in particular regarding attention to display ads by Huberman et al. (1998) argues that online consumers have a lower threshold for uncertainty early on in a web browsing session, and thus advertisers should expect more clicks at the beginning of a session. But this result is in contrast with later work by Chatterjee et al. (2003) who find clicks to be equally likely throughout the browsing session, indicating that timing of ads throughout a browsing session does not matter for ad effectiveness.

In addition, Chatterjee et al. (2003) point out that display ad click-through rates are quite low and they have been dropping over time. This claim increases the need for us to better understand online advertising engagement. The falling engagement rates of display ads have led many advertisers to more obtrusively interrupt consumers with pop-up ads which co-opt consumer attention, often against their desires. Moe (2006) studies the effects of pop-up ads on engagement and show that they initially improve customer responsiveness to display ads but also stop working over time.

This type of obtrusive interruption of web browsing is illustrative of the conflict between the goals of consumers and those of advertisers. Danaher and Mullarkey (2003) investigate more thoroughly the goal-directed nature of consumers' online browsing behavior. They find that the more goal-oriented the browsing session, the less display ads will be attended to. Similarly, Cho and Cheon (2004) find that consumers avoid attending to display ads because it impedes their browsing goals. Collectively, this stream of research points to the idea that, in order for display ads to get initial attention, consumers often need to (deliberately or not) momentarily give up the attention they allocate to their primary task of choice that led them to browse the internet.

1.3 Attention to Tasks

We now turn to examine attention and how it relates to everyday tasks that are achieved through various types of media such as print, television and the Internet. We begin by defining three concepts to aid in understanding task behavior. First, we adopt the definition of a *task* as a complex situation capable of eliciting a goal-directed performance from an individual (Fleishman and Quaintance 1984). Second, we define *task progression* as a collection of steps that make up the task, all of which need to be completed in order to achieve the task's goal (Cooper and Shallice 2000). Task progression allows us to label the first step, any middle steps and the last step in a task. Finally, we define *task type*, which allows us to recognize that within any period of activity there may be multiple different tasks, each with its own progression.¹ We illustrate the application of these constructs in detail in the measurement section to follow.

With these definitions in mind, we review the growing body of literature in psychology around understanding the mechanism of attention, the effects of external stimuli in capturing attention and the task engagement mechanism.

In this literature, attention is often discussed in terms of perceptual load theory, whereby individuals are endowed with a fixed capacity for attending to tasks. Each task has an attentional cost, and so attention to external stimuli depends on available attentional capacity (Lavie, 1995; Lavie and Tsal, 1994). The more an individual focuses on a task, the more attentional resources are allocated to that task, and the less likely they are to have excess attentional capacity to attend to other stimuli (Forster and Lavie, 2009).

¹ Classification of tasks is particular to the environment being studied and Fleishman and Quaintance (1984) discuss in detail the classification of tasks based on four dimensions: (1) task content, (2) task environment, (3) level of learning, and (4) discriminable task functions.

In this model of attention, even if there is spare attentional capacity, any external stimuli will almost always be ignored if it can be easily identified as not relevant to the task at hand (Lamy et al., 2004). At first, this process is deliberate. However, eventually, individuals will adopt habits to automatically ignore task-irrelevant stimuli (Anderson et al., 2011).

The fixed-attentional capacity model can help explain online ad engagement, or lack thereof. For instance, the banner blindness findings of Benway and Lane (1998), goal-orientation findings of Danaher and Mullarkey (2003), and why people ignore ads while engaged in a task are predictable outcomes of the fixed-attentional capacity model. Banner blindness is a behavior by which consumers who are attending to an online task (e.g., reading text) end up suppressing attention to banner ads. Similarly, the goal orientation findings are based on a treatment condition of a highly engaging task versus a control condition of no specific task other than casually browsing web pages.

In a closely related paper, Tavassoli et al. (1995) investigate the ability to process information contingent upon television program involvement. They find an inverted-U relationship between ad recall and program involvement. Their research is similar to ours in that they measure how memory and attitude towards TV advertisements varies with primary task (e.g., TV viewing) involvement. However, they examine attention *across* consumers, whereas we investigate attention *within* consumers.

In the Norman-Shallice model of executive control (Norman and Shallice 1980, 1986), action statements (i.e., schemas or tasks) are selected by the individual's executive control system. Once selected, a task remains active until its goal is reached or it is inhibited by a competing task. We contend that online tasks such as reading emails, browsing news articles, and watching videos are deliberately selected by Internet users that go online. Rarely, if ever, does a person have as

their primary task to pay attention to banner ads. These are competing tasks that can temporarily inhibit a primary task.

Attention allocation to primary tasks slowly builds up as consumers engage with them (activation in the Norman-Shallice model) and stay high until the primary task is completed, abandoned, or attention is co-opted away with competing stimuli (deactivation in the Norman-Shallice model). Subsequently, Cooper and Shallice applied the Norman-Shallice model to explain cognitive resources build up and release in day-to-day tasks such as coffee preparation (Cooper and Shallice 2000, p. 319). We contend that the same general pattern applies to online tasks as well. If so, a natural question to ask is how does online task progression impact engagement with display ads shown in the beginning, middle and end of tasks commonly done online?

1.4 Conceptual Model

Given the above findings regarding task and attention, we now propose to relate them to online browsing behavior through a simple and generalized conceptual model. Since an online task is accomplished via attending to an ordered sequence of web pages (e.g. reading several news articles sequentially on CNN.com), the most obvious choice for measuring task step is web page visit, where each distinct page view is a unique task step. Task progression measurement entails that the first page of the task be labeled a first step, the last page view of the task be labeled a last step, and all other page views in-between be labeled as middle steps.

General attention theory states that until a task is completed, the consumer is invested in and attending to that task. Throughout the online task, ad attention is low while the consumer is actively engaged in their primary task. Since attentional resources revolve around a task and not

the browsing session or the domain, breaks between task steps (such as moving between web pages or domains) as part of a single task do not have a significant effect on ad attention, all else equal.

Eventually, as people complete one task, they either move onto the next task, or they spend some (free) time wandering around until they perceive a need to engage in another task. In both situations, they are “released” from the task and their attentional resources free up. At this point, people disengage from their primary task and top-down processes (i.e., individual-dominated attention) will again take-over in choosing the next task. Collectively, these corollaries to the theory drive our central hypothesis. We predict that, within a task, attention to ads will be higher at the first and last task steps than at any other task step in between. The reasoning is that it is at these transition points between tasks where people have not fully assigned their available attentional capacity to a primary task of choice, thus leaving unallocated attention for other stimuli, ads being only one of them. We note that our hypothesis pertains to relative attention within a task and does not speak to levels of attention across tasks.

In reality, many ads appear when consumers do not have the attentional availability to engage with them. By placing ads in the middle of a task, advertisers are attempting to capture a consumer’s attention when they are most engaged in another task and thus least likely to attend to the ad. This conjecture would help explain the high ad avoidance in browsing sessions described by Cho and Cheon (2004) as most ads occur mid-task. A challenge in testing this hypothesis is how to classify tasks and determine their beginning and ending points. We will address this concern separately for each of our two data sets.

1.5 Measurement

To test our central hypothesis, we operationalize the task progress construct as follows. We define a user-session-page observation as the unit of analysis. A session is considered a sequence of web page visits where the time between page visits is less than 5 minutes.² Presumably, people will engage in many browsing sessions throughout the day. Within each browsing session, people may have multiple tasks. We assume that people engage in and complete one task at a time. We acknowledge that many people engage in online multitasking, where they simultaneously view pages from different tasks or domains. In our data this is somewhat rare³ and so we exclude them from our sample for the sake of parsimony.

Not being able to ask people about the tasks they are performing online, we determine the type of sequential tasks based on the content category of each website. This is directly in line with the Cooper and Shallice (2000) model describing the main task as the ‘goal’ that people wish to achieve. For example, a person may choose to read news websites so we assume the task goal to be ‘being informed about current events,’ an online gaming site allows accomplishing the task of ‘entertaining oneself by playing online games,’ a weather site allows ‘knowing the weather.’ Similarly, the Cooper and Shallice (2000) model implies that any other page visit within the same content category is a subtask and considered a step towards completing the overall task.⁴

² We compared break intervals from 2 to 10 minutes and observed no qualitative differences.

³ 100% of the sessions in study 1 and 98.7% of sessions in study 2 contain a single task at a time. This is a conservative measure of multi-tasking as sessions with serially completed tasks are a part of the 1.3%.

⁴ We exclude so-called online portals that let people accomplish multiple different tasks due to task identification concerns.

To illustrate task progression, we decompose a hypothetical web browsing session presented in Figure 1-1. In this example, session 1 shows a user who starts a browsing session with the goal of being informed about current events; the resulting task is ‘News.’ They then proceed to visit the web domain CNN where step 1 is the first Uniform Resource Locator (URL) page that they visit. This step also corresponds with the first step in task progression. In step 2, the person remains on CNN and navigates to a new URL page. Step 2 is still part of the ‘News’ task and a middle step in terms of task progression. Next, the user navigates to the BBC domain by loading a new URL in step 3. This again remains the task of ‘News’ and step 3 is also a middle step in terms of task progression since the task has not changed (even though the domain has changed). Step 4 is a new URL page on the BBC domain. Being the final ‘News’ step, it is classified as a last step in terms of task progression. Within this same session, the user begins the new goal of playing online games. They navigate to puzzles.com and the first URL page loaded is labeled step 1 of a ‘Games’ task which is again the first step in terms of task progression. From here the labeling process carries on. It should be clear that first and last steps in terms of task progression are distinct from domains and web session first and last web pages.

For the case in Figure 1-1, our hypothesis predicts that within each task the shaded cells (e.g., Session 1, News at steps 1 and 4) have higher levels of ad attention than the unshaded cells. Note that this is distinct from the industry belief that ads at the top and bottom of web pages have higher attention levels (Brebion 2018, Work and Hayes 2018) and from prior research showing that ad engagement is equally likely throughout the browsing session (Chaterjee et al. 2003).

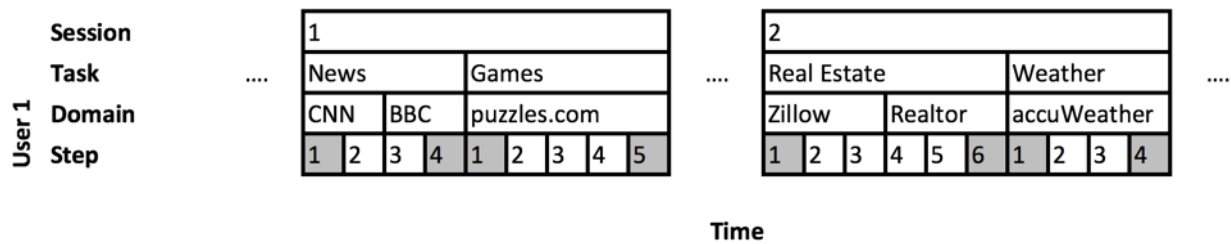


Figure 1-1 Task Progression Conceptualization

1.6 Design and Identification

Digital advertising impression logs are a natural data source for this type of investigation. First, they allow the researcher to track users across web pages and time. Second, web page content provides natural task classifications. And third, digital ads are external stimuli within each step designed to grab a user’s attention. With this type of data, we construct the user browsing session, classify tasks within that session and identify task steps as defined above to test our hypothesis.

To measure ad engagement we employ two metrics: clicks and video view-throughs. Clicks have been a standard metric in marketing research and have the benefits of being a behavioral measure (Chaterjee et al. 2003). More importantly, clicks are a direct measure of engagement with each ad. Since our theory makes predictions for ad engagement at each task step, our measure needs to be fine enough to capture changes in ad engagement at the granularity of a task step. Clicks allow for this.

Clicks have become more contested lately with the increased focus on purchase attribution (Hoban and Bucklin 2015). Given the nature of the web, it is very challenging to attribute a purchase to a single ad, and some have argued that purchase events should be attributed to multiple

ads (Li et al. 2016). This makes purchases a poor measure for our purposes as our hypothesis concerns engagement with individual ads.

For video ad formats, we propose a second measure, ad view-through. This is a binary measure of whether or not a user finished viewing a video ad. It has gained popularity in industry as an alternative to clicks for videos (Heine 2014). This measure is also closely related to TV ads and fits naturally with the traditional marketing literature. Another advantage of measuring view-throughs is that they occur more often than clicks and therefore allow better measurement of advertising engagement on smaller datasets. Lastly, using two measures allows us to test our hypothesis on both an opt-in metric (clicks) and an opt-out metric (auto play view-throughs).

We test our hypothesis with two data sets representing different observational designs. The details of these datasets are discussed in the subsequent sections, but we briefly mention the design intentions. The first data set was collected as part of a separate ad server experiment. Thus, it is small, without missing data and the ads were assigned in a controlled and randomized manner. However, it is restricted to a single task (news reading) and single domain per session resulting in the beginning and endings of the browsing session, task and domain effects being confounded. This design represents a simplified situation in a controlled environment, and is represented in Figure 1-2.

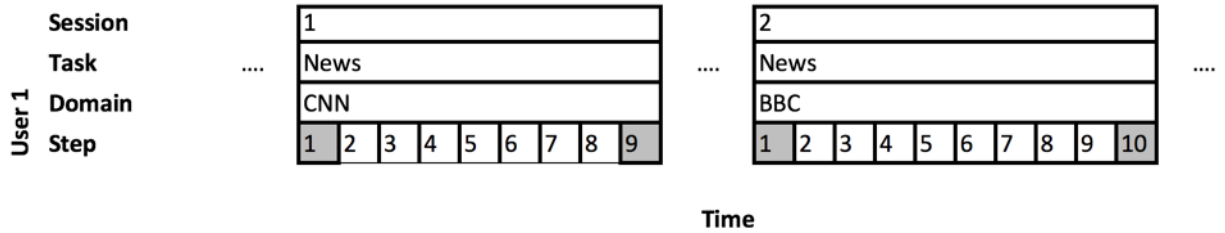


Figure 1-2 Data Set 1 Task Progression Identification

The second data set reverses these strengths and weaknesses. This data set is much larger but misses data due to the ad server operating in a real-time bidding (RTB) environment. A larger sample size allows us to examine multiple task types and to tease out the differences between the beginning and endings of session, task and domain effects. This observational study design is represented in Figure 1-3, where the shaded cells represent domain and session step. Doing so allows for the identification of the task step effect separately from potential domain and session confounds.

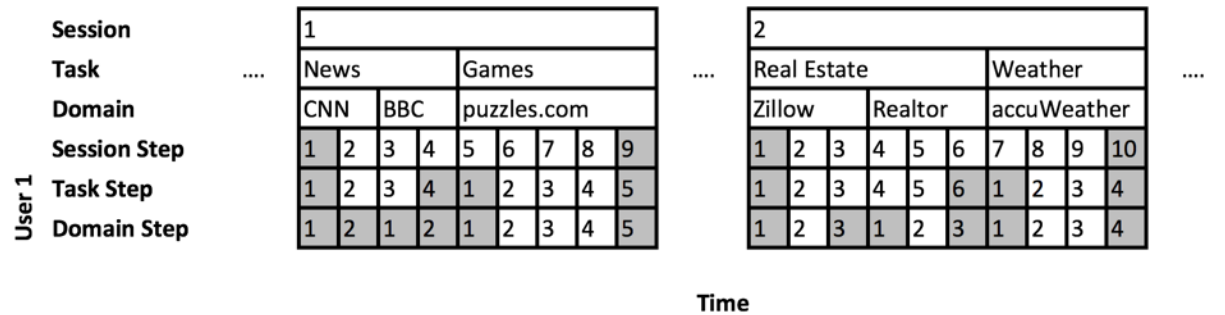


Figure 1-3 Data Set 2 Task Progression Identification

Given that the above design identifies the treatment effects of interest, our next concern to address is whether or not the observations in each dataset are truly random. Traditionally, when modeling user-level behavior with clickstream data, it is assumed that observations are

independent (Chaterjee et al. 2003, HobanBucklin 2015). But this is not necessarily true for all online browsing behavior. Web pages and their ads are not randomly assigned during a browsing session. Websites often have hierarchical structures beginning with homepages possessing generic content and branching out deeper with specialized content. Users tend to enter websites near the top of these "tree-like" structures, generally at a homepage.⁵ Following different routes, users then continue deeper before finally leaving the website from one of its branches. Consequently, observations pertaining to attention to ads on web pages are not in general independently and identically distributed.

As an example, consider a typical visit to the website dictionary.com (contained in our second dataset). A visit to this site often starts at the homepage where a search query is entered, then a results page is loaded with a list of possible words. This is typically followed by the user navigating to one or several of these word definition pages. From this example, we see that a potential confound in modeling engagement to ads is that the types of web pages that are likely to be first steps, middle steps, and last steps may be structurally different from each other in terms of content and, more importantly, ad type, load, and placement. As a consequence, any constructed dummy variable for measuring task progression is also going to unintentionally capture systematic page structure effects.

We address this confound by applying the reasoning put forth by Rubin (1974) in propensity score matching. The premise is that we want our dataset to resemble a randomized control trial where one can stratify across all control variables for treatment, in our case an ad shown during a certain task progression (e.g., beginning) and a control group, in our case that same

⁵ This is not true for all websites, but a large proportion of websites tend to be structured in this tree-like fashion.

ad shown at a different point in the task progression (e.g., middle). This procedure creates a balance across all control variables which is the foundation of a random assignment that allows us to estimate the causal effect of the treatment. The result is a data set that contains probabilistically identical samples between treatment and control groups. Another alternative is to control for balancing variables (i.e., page structure) in the model specification directly. However, Dehejia and Wahba (1999) illustrate that merely controlling for such variables in the model specification is not sufficient to adequately account for the bias in observational samples.

A final concern is that of endogeneity at the end of a task. Is it, in fact, the ad engagement (e.g., view-through or click) that caused the end of the task rather than the end of the task causing more engagement? This concern is mitigated by the propensity scoring matching procedure. By building a model to predict the beginnings and ends of tasks based on the domain and URL structure, and not any measures of ad engagement, we test the effect of ad engagement on task progression. Intuitively, if task progression were unimportant, then we would not be able to predict which pages would likely be the beginnings and ends of tasks. However, if one can predict task progression, then it implies that it is page structure and not ad engagement that predominantly drives where consumers enter and exit websites. Additionally, observations where ad engagement did cause the end of a task are more likely to be excluded in the matched datasets since they are less likely to be predicted first/last steps.⁶

⁶ The authors would like to thank Donald Rubin for helpful comments and suggestions regarding propensity score matching and identification.

1.7 Data Set 1 – Single News Tasks

Our first data set is made up of proprietary video advertising data from Vuble⁷. This dataset provides us with the full browsing history of users for several news websites, described in more detail in the following section. While this limits us to a single task, it also allows us to precisely map the sequence of distinct URLs requested by each user to task steps and thus fully identify task progression for each user in each session. As a limitation of this dataset, an entire URL page view (not the top or bottom of a page) is defined as our task step since we do not have tracking data for users as they move within a single page.

Vuble uses a variety of video ad formats, but pertinent to our analysis are the ones where a video ad was placed at the top of the page. As users scroll down a web page it pauses and disappears. Then, it reappears and restarts once the user scrolls to the bottom of the page. This feature allows us to measure how long a user spends at the beginning and end of a web page. Had the ad only played at the top of the web page, as is common, we would not be able to measure the time spent engaging with an ad at the bottom of the web page.⁸

1.7.1 Data Description

The data was collected over a period of seven days, from April 15 to 21 of 2016. It consists of 10 different video ads of various consumer brands ranging in length from 25 to 32 seconds long collected from 5 news websites in the USA and France. All of these videos have the same viewable

⁷ Previously known as Mediabong.

⁸ In the case of reading the news online, as a user gets to the last step of a task, they scroll down and continue to read the article. At the start of the last step, they are still highly engaged in news reading and are expected to scroll past the ad which is then paused. By continuing to play the video ad at the bottom of the page the advertiser can re-engage the user at the time they expect them to leave the task and the ad would have higher attentional capture. If a video ad did not continue at the bottom of a web page at the last step, where many people would end a task, then our measurement of engagement for the last step would be understated.

size of 500 x 280 pixels. The ads were shown in the morning (10 am to 11 am), afternoon (2 pm to 3 pm) and evening (6 pm to 7 pm) to control for outside attentional demands. These ads were shown across time of day, day of week, country and campaign in a randomized and counterbalanced manner to control for content and time of day effects. Users were unaware of any specific advertising study as it was implemented on the publisher's web pages and shown to users at random.

We collect all data at the user-ad impression level. In total, we observed 63,402 browsing sessions consisting of 71,508 page visits and 55,561 uniquely identified users. Of these sessions, 33,766 consisted of a single page visit and thus cannot be used. Additionally, 23,688 of the page visits were of the type where the video ad restarted at the bottom of the page. Since our analysis depends on at least a first, middle and last step as well as engagement measures at the top and bottom of the page, we reduce the data set to 4,527 page visits from 3,011 users and 1,246 unique news articles. However, we still use the full data set to construct the browsing sessions and identify the first and last task steps. As can be seen here, our empirical analysis requires a very large dataset size to apply the propensity score matching portion of the model.

We define an end to a browsing session after five minutes of user inactivity and note that the main results of interest are robust to varying choices of inactivity intervals commonly defined in the range of 2 and 10 minutes. As our measure of advertising engagement, we consider video ad view-through rate (VTR) which is defined as the percentage of video ads that completed playing before the user navigated away from the web page. We also compared other ad engagement metrics (percentage of video watched, time watched, paused, muted, etc.) and note that the results are qualitatively similar.

Our hypothesis is that ad engagement should be higher at first and last task steps relative to a middle step. In a cursory test of this hypothesis, we classify task progression based on the observed sessions and calculate the mean video ad completion rate for the first, middle and last task steps. We observe that the video completion rate of the first step is 38%, which drops to 22% for the middle step and then climbs again to 34% for the last step. The U-shape of these point estimates provide initial support for our hypothesis, but we note that the large standard errors (ranging from 41% to 49%) around these estimates result in no statistically significant differences between the three groups. Therefore, we require a model to control for other known sources of variation such as ad, user and website specific effects.

For completeness, we note that sessions consisting of a single page visit have a video completion rate of 38% with a standard deviation of 48% indicating a similar response to that of a first task step. This provides some initial indication that targeting users at the beginning of task may be beneficial for an advertiser regardless of session length.

1.7.2 Empirical Model

Next, we propose a model to estimate user engagement throughout tasks. We begin by recalling our hypothesis that users attention starts low, builds up and is eventually released again at task completion. Therefore, we predict that user engagement with ads will behave in a U-shape: higher at the beginning and ends of online tasks relative to any step in the middle. We now proceed by specifying an empirical model to test this hypothesis by measuring ad engagement contingent upon online task progression.

We measure online ad engagement at the web page level by following users across browsing sessions. Our dependent variable is the binary outcome of whether or not the video ad completed playing (VTR). Our unit of analysis is the user-session-webpage combination for a

single task. For each user i in browsing session j we have a time-ordered sequence of web page views k . We look up the content category for each web page to assign a task type and only consider one task at a time (i.e., the no multi-tasking assumption). The first (last) web page viewed in each task for each user is defined as the first (last) step in the task. Throughout the paper we use dummy variables to code for first and last task steps. Their coefficients in the regression directly relate to our hypothesis and are the main parameters of interest.

We model online user behavior using logistic regression with random effects, which has now become standard in the clickstream literature (Hoban and Bucklin 2015). In examining the data set for sources of VTR variation, we find notable individual campaign, campaign type, and domain effects. We note minor variation in time of day and day of week effects and control for them as well.

In line with past research and our data exploration, we include the standard fixed effects for website domain (e.g., *cnn.com*), time of day and day of week variables, as well as campaign type variables (ad and brand features). We then include random effects components for user ID, campaign ID and URL to capture the main individual heterogeneity effects. We model ad engagement as:

$$P(\text{Engagement}_{ijk} = 1) = \alpha + \beta_1 \text{Website}_{ijk} + \beta_2 \text{Time}_{ijk} + \beta_3 \text{CampaignType}_{ijk} \\ + \gamma_1 \text{FirstStep}_{ijk} + \gamma_2 \text{LastStep}_{ijk} + \delta_1 \mu_i + \delta_2 \mu_{\text{campaign}} + \delta_3 \mu_{\text{URL}} + \varepsilon_{ijk},$$

where α is the intercept, β are fixed effect parameters, γ are the treatment effects being tested, δ are random effects parameters with $\mu \sim N(0, \Sigma)$ and ε is the standard logistic idiosyncratic error term.

1.7.3 Results

The output of various logit models is shown in Table 1-1. We estimated model 1, with domain fixed effects only, model 2 with domain, day of week and day-part fixed effects, and model 3 with domain and campaign type fixed effects. Model 4 is the full random effects model. Model 4 appears to be the best fit based on the AIC scores, while model 2 has very close AIC scores as well. In line with our hypothesis, both first and last step coefficients show significantly higher rates of video completion than medium steps. These results are robust to all specifications estimated. The results of model 4, in particular, indicate an 82% increase in the probability of completing a video ad at the beginning of tasks and 30% increase at the end of a task when compared to a middle step. In terms of the covariates used, in support of previous findings (Chatterjee 2003), we report no significant time effects. As expected, we do find that some campaign types appear to be more effective than others for capturing attention.

These results provide initial support for our hypothesis that consumers are more likely to pay attention at the beginning and ends of the online task when it concerns the task of reading news online. These results also appear to be independent of other important variables that are known to impact online ad effectiveness. Unfortunately, we cannot claim causality as these estimates do not arise from a randomized sample of first, middle and last steps. In the next section, we address this limitation by applying propensity score matching to the Vuble dataset.

Table 1-1 Data Set 1 Model Output

	<i>Dependent variable:</i>			
	Video Completed			
		<i>logistic</i>		<i>generalized linear mixed-effects</i>
		(1)	(2)	(3)
First step	0.592***	0.589***	0.595***	0.602***
Last step	0.268**	0.266**	0.259**	0.265**
www.economiematin.fr	1.358***	1.359***	-0.011	0.002
www.independent.co.uk	0.512**	0.516**	-0.452	-0.446
www.livefoot.fr	-0.189	-0.187	-1.566**	-1.553**
www.remedes-de-grand-mere.com	0.720***	0.712***	-0.670	-0.663
Saturday		0.196		0.175
Sunday		-0.025		-0.049
Monday		0.108		0.051
Tuesday		-0.024		-0.036
Wednesday		0.033		0.003
Thursday		0.106		0.063
Midday		-0.077		-0.070
Morning		-0.049		-0.042
Campaign type: Beauty&Health			1.539**	1.529**
Campaign type: Car			1.326*	1.322*
Campaign type: CPG			0.827	0.832
Campaign type: Electronics			0.729	0.742
Campaign type: Entertainment			1.232*	1.231
Campaign type: Supermarket			1.180	1.166
Campaign type: Telecom			1.619**	1.613**
Constant	-1.538***	-1.552***	-1.536***	-1.542***
Observations	1,956	1,956	1,956	1,956
Akaike Inf. Crit.	2,159.746	2,173.459	2,162.288	2,175.011
Bayesian Inf. Crit.	2			2,303.320

Note:

*p<0.1; **p<0.05; ***p<0.01

1.7.4 Propensity Score Matching Design

In our application, we can see task progression as our three treatment categories: first, middle and last steps. By setting the middle step to be the baseline (control) we choose first and last steps to be two distinct treatment types. Next, we need to create two new balanced data sets, one for each treatment. First, we model $P(\text{treatment}=\text{first})$ and build a matched dataset of middle and first steps. We then repeat this process with $P(\text{treatment}=\text{last})$ and build a matched dataset of middle and last steps. These two balanced datasets would provide us with a quasi-experimental design setup for making causal inferences about the first and last step treatment effects (Rubin 1974).

Using the insight that users enter and exit websites according to their structure, we are able to partially predict entry and exit using information from web pages' URL string. For example, `www.dictionary.com` is probabilistically more likely to be a first step; whereas `www.dictionary.com/browse/quandary/page=synonym` is more likely to be a last step. We construct variables that tend to explain where people are entering and exiting a website.

We code three variables associated with domain structure characteristics. We begin by defining a *page depth* variable as the count of forward slashes in the URL. A second variable that we code are *page count* indicators. These are string components of the form 'page=' or 'p='. These are also associated with someone moving deeper through a website, which is generally not present in a first step URL. Our third constructed variable is the *step number* of the given session. When we identify our browsing session and all the pages it contains, we can then order all the pages as discrete page steps navigated across time. By simply numbering these pages in ascending order as they occur in a session, we get a step number. We can then use this number as a discrete timestamp to predict the likelihood of ending the session. The higher the ordered step number of a web page,

the more likely is the user to have exited a website on that page. This is operationalized by including a polynomial of step number in $P(\text{treatment}=\text{last})$ such that we get a hazard function over time as a session progresses.

We model the last step as a function of task domain, step number, page depth and page count. We omit step number in the first step model as a step count equal to 1 is equivalent to our dependent variable in the propensity score model. Since the purpose of propensity score matching is maximizing prediction performance (i.e., in-sample fit) and not deriving interpretable parameters, we implement a random forest algorithm⁹. Before creating the final matched datasets we linearize the propensity scores using a logit function to make them approximately normally distributed. We use these models to predict the probability (propensity) of each observation of being a first versus middle step and last versus middle step.

1.7.5 Results of Propensity Score Matching

We now turn to building a data set using propensity score matching. As discussed above, by building a propensity scoring model, we endeavor to estimate the underlying assignment mechanism as it may lead to dissimilar observations between treatment groups. One way to validate that the model described above captures this mechanism is to test our prediction accuracy out of sample. A model capturing the underlying assignment mechanism should be able to correctly classify task progression in a holdout set with a high degree of accuracy. Following standard practice, we partition our data into 80% training and 20% holdout sets (Hastie et al. 2009). We then fit the random forests algorithm on the first step of the training set using the page features

⁹ The random forest algorithm (Breiman 2001) creates an ensemble of decision trees using bagging to construct a set of trees with controlled variance. This algorithm samples both observations and features using bootstrapping and then averages the prediction results. This method is particularly robust to overfitting and tends to perform very well relative to other parametric models at prediction tasks (Segal 2004).

discussed previously. The resulting receiver operating characteristic curves for the holdout sets are shown in Figure 1-4. Applying this model on the first step test set yielded a 0.73 area under the curve (AUC) score¹⁰. We do the same for the last step, this time adding step count as an additional variable. This yielded an AUC score of 0.78. These high AUC scores provide strong support for our assumption that the URL structure variables are strong predictors of task progression. These out of sample results also illustrate the ability to predict task progression in practice.

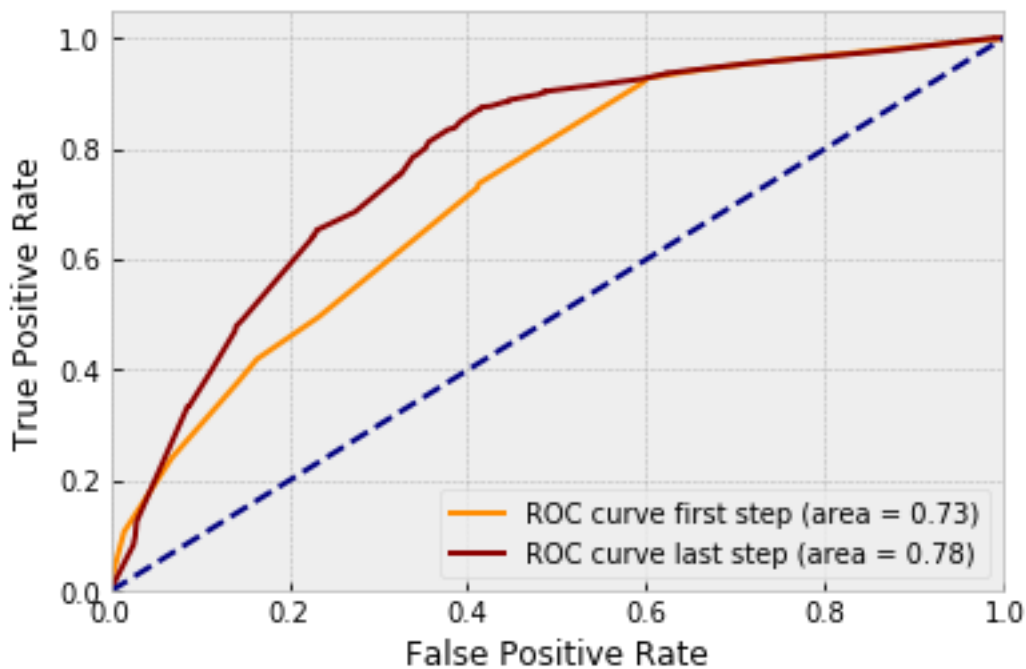


Figure 1-4 Receiver-Operator-Characteristic Curves for Task Progression

With these appropriately matched data sets, we are now able to repeat our previous analysis, and can now more accurately show how task progression influences the attention to

¹⁰ The receiver operator characteristic (ROC) curve plots the true positive rate vs. false positive rate for the prediction of a random variable. The area under this curve is used as a measure of predictive accuracy where the random guesses give an expected AUC of 0.5, and perfect prediction gives an AUC of 1.

display ads. We again estimate the four model specifications on each data set. The results are shown in Tables 1-2 and 1-3 for the first and last steps respectively.

Table 1-2 shows the results for the first step matched data where middle step is the baseline. We can now interpret the first step parameter as the treatment intervention and regard it as the effect of targeting consumers at the beginning of a news-reading task on their attention to the ad. We first note that this parameter is significant in all four model specifications. The odds of completing a video ad in the first step are now 58% higher when compared to the middle step. Recall that it was estimated to be 82% higher in the unbalanced results in Table 1-1. This difference in magnitude can be attributable to the effect of the web page confound that propensity scoring has controlled for, *dampening* the non-causal estimated impact. There is indeed a strong web page level effect on attention to ads for the type of web pages that users are likely to start a task on.

Table 1-3 shows the analogous results for the last step matched data. We similarly note that the last step parameter is significant under all four model specifications. We again see a markedly different magnitude of coefficients due to the matched data. The odds of users completing a video ad in the last step are now 50% higher when compared to the middle step instead of the 30% higher in the unbalanced results in Table 1-1. The Rubin-style causal inference procedure de-biases the result in this case by *enhancing* the magnitude of the non-causal estimation, as opposed to dampening it. The reason is that typical last step pages tend to have lower ad engagement rates relative to typical middle step pages. When we use propensity score matching, we create a reduced data set where each treated observation has a corresponding control observation that has the same expected engagement rate given the page structure. This removes the bias of the task progression parameter since it is no longer capturing page effects. For both matched data sets, a power test was conducted and showed that the power to detect effects of this size exceeds 99%.

Substantively, it is also noteworthy that the increased odds of completing online video ads are now roughly the same size for both the first and last steps once controlling for page level confounds. This suggests that it is not that first or last steps evoke higher attention to ads than middle steps but rather the opposite. Middle steps, when a person is deeply engaged with a task, are when attention allocated to ads are diminished. As such, this result strictly speaking does not help advertisers determine when during a consumer's task to advertise; but rather it informs when they should *not* advertise. While precisely controlled, the limitations of this finding, as it stands, are that it only applies to the task of reading news online in very simple browsing sessions and that session and domain effects are confounded with our task effect. In the next section, we address these shortcomings by incorporating several of the most common tasks that consumers perform online in more complex browsing sessions containing multiple domains, disentangling their effects.

Table 1-2 Data Set 1 Propensity Score Matching Model Output for First Steps

	<i>Dependent variable:</i>			
	Video Completed			
		<i>logistic</i>		<i>generalized linear mixed-effects</i>
	(1)	(2)	(3)	(4)
First step	0.455***	0.452***	0.459***	0.456***
www.economiematin.fr	1.143***	1.144***	-1.247	-1.206
www.independent.co.uk	0.449	0.453	-1.240	-1.221
www.livefoot.fr	-0.384	-0.382	-2.799**	-2.754**
www.remedes-de-grand-mere.com	0.638***	0.635***	-1.756	-1.714
Saturday		0.014		-0.046
Sunday		0.246		0.207
Monday		0.377		0.297
Tuesday		-0.053		-0.093
Wednesday		0.093		0.078
Thursday		0.067		-0.012
Midday		-0.211		-0.198
Morning		-0.045		-0.035
Campaign type: Beauty&Health			2.569**	2.526**
Campaign type: Car			2.590**	2.558**
Campaign type: CPG			1.367	1.348
Campaign type: Electronics			1.109	1.085
Campaign type: Entertainment			2.156*	2.120*
Campaign type: Supermarket			2.222*	2.164*
Campaign type: Telecom			2.484**	2.426**
Constant	-1.332***	-1.343***	-1.334***	-1.312***
Observations	892	892	892	892
Akaike Inf. Crit.	1,094.664	1,107.610	1,094.590	1,109.950
Bayesian Inf. Crit.				1,215.406

3

Note:

*p<0.1; **p<0.05; ***p<0.01

Table 1-3 Data Set 1 Propensity Score Matching Model Output for Last Steps

	<i>Dependent variable:</i>			
	Video Completed			
		<i>logistic</i>		<i>generalized linear mixed-effects</i>
	(1)	(2)	(3)	(4)
Last step	0.401**	0.387**	0.418**	0.407**
www.economiematin.fr	1.170***	1.147***	0.477	0.305
www.independent.co.uk	0.395	0.394	-0.150	-0.230
www.livefoot.fr	-0.290	-0.320	-1.018	-1.202
www.remedes-de-grand-mere.com	0.382	0.392*	-0.379	-0.511
Saturday		0.261		0.250
Sunday		-0.534		-0.588
Monday		-0.566		-0.614
Tuesday		0.051		0.085
Wednesday		-0.080		-0.125
Thursday		0.005		-0.005
Midday		0.147		0.168
Morning		-0.083		-0.081
Campaign type: Beauty&Health			0.714	0.858
Campaign type: Car			0.159	0.284
Campaign type: CPG			0.866	0.823
Campaign type: Electronics			1.160	1.308
Campaign type: Entertainment			0.751	0.902
Campaign type: Supermarket			0.484	0.643
Campaign type: Telecom			1.358	1.530
Constant	-1.530***	-1.463***	-1.540***	-1.465***
Observations	916	916	916	916
Akaike Inf. Crit.	970.735	978.755	968.703	977.605
Bayesian Inf. Crit.				1,083.645

4

Note:

*p<0.1; **p<0.05; ***p<0.01

1.8 Data Set 2 – Browsing Sessions with Multiple Tasks

This second data set comes from MaxPoint, a digital advertising platform that buys online ads on real-time bidding (RTB) exchanges on behalf of its client advertisers. The primary benefit of data from ad exchanges is that they offer ads from millions of websites which allows us to track users across browsing sessions, as they move between websites. The downside is that, unlike data set 1, the ad server does not guarantee full transparency into the entire browsing session. We only observe those instances in which MaxPoint won the right to place an ad. Since many companies compete for each ad opportunity on these exchanges, MaxPoint will only observe ad engagement data for ads that it won on the auction. If this data were missing completely at random (MCaR) then this would not bias any results, only adding more noise, which can be overcome by increasing the sample size (Agresti 2015). If it is not MCaR, then the parameter estimates could be biased.

We make the argument that this data is missing completely at random for the following reason. In talks with senior managers at MaxPoint, they have no reason to believe that their bidding algorithms favor more, or less, people at the beginning, middle or ending of user tasks. In fact, they claim to not favor any of these ad placement opportunities over others, to the best of their knowledge. Neither do they track, measure or decide based on proxies correlated with task progression (more on targeting in the next sub-section).

We also note that missing data cannot result in first or last task steps being misclassified as middle steps. This claim is derived from the facts that a task is an ordered sequence of pages and that the treatment effects are themselves an ordered sequence. Thus, if we miss a first step in a task, it will result in a middle step page to be called a first step page. If a last step page is missed

in a task, then it will result in a middle step page to be called a last step page. Finally, if a middle step page is missed in a task then the other first, middle, and last step pages are still correctly classified. With middle step as the baseline in all our models, any missing observations will only lead to the treatment effect estimates to be understated.

1.8.1 Data description

The data was collected over a period of fourteen days, from November 22nd through December 5th in 2016. The ads were shown across the entire 2-week period across desktops, tablets and smartphones, and observed through the regular course of business. The ad sizes are the industry standard of 160 x 600, 300 x 250 and 728 x 90 and we observe RTB exchange data indicating whether the ad placement was above or below the fold, or unknown.

Specific to this data set, MaxPoint uses two types of targeting strategies: user targeting, which delivers ads contingent upon specific user characteristics, and a more general brand awareness targeting strategy that does not rely on user or page characteristics. For modeling purposes, we exclude all observations arising from user profile targeting. Additionally, it is often the case that multiple ad campaigns desire the same ad placement. In this case, the tie is broken using random assignment. We also note that there are no statistically significant price differences related to session, task and domain progression.

For this data set, we use clicks as the variable of interest for measuring ad engagement since we have enough observations to reliably estimate click-through rates.¹¹ For this analysis we

¹¹ We note that a part of MaxPoint's platform is a proprietary algorithm to detect non-human traffic. We used this information to remove non-human observations.

selected 20 website categories¹² to define our online tasks. We chose these particular tasks so as to cover a wide variety of online behaviors (such as browsing dating sites, real estate shopping, playing online games, etc.) and verify that our theory generalizes well beyond just news reading.

We collect all data at the user-ad impression level. In total, we observed 105,846,160 page visits from 8,066,795 uniquely identified users and 839 advertising campaigns. We only consider display ads collected from 20,928 websites in the USA representing 20 tasks that were judged to map content of the websites. We are again only interested in sessions with at least 2 page views (which include distinct first and last task steps) and sessions with at least one of our 20 tasks.

Using the above criteria, we reduce the data set to 34,116,671 page visits. However, we again use the full data set to construct the browsing sessions and identify the first and last task steps. We again define an end of a browsing session as five minutes of user inactivity.

Table 1-4 below summarizes the session data for these tasks. We see a substantial amount of variation in activity across tasks.

¹² Website category information was provided by Amobee (a digital marketing intelligence company and ad server), formerly Turn.

Table 1-4 Data Set 2 Summary Statistics

Task	Domains	URLs	Users	Campaigns	Impressions	clicks
Automotive	906	85,955	112,987	726	817,792	818
Careers	161	18,657	20,262	509	140,299	40
Dating	59	13,248	5,410	340	78,755	44
Education	1,173	138,021	161,747	815	1,204,347	2,580
Entertainment	4,697	655,201	1,138,318	896	8,154,216	17,122
Food	1,618	156,323	258,518	847	1,297,011	5,205
Games	1,100	146,326	89,172	800	947,134	1,687
Health	1,014	81,582	129,300	792	799,174	3,792
Hobbies	778	104,676	205,729	855	2,167,323	2,489
Home & Garden	327	31,552	40,158	597	186,389	43
Music	366	20,916	53,327	705	589,280	699
News	3,752	624,970	1,510,322	896	10,525,330	19,369
Pets	281	15,773	31,028	568	125,360	51
Real Estate	163	143,964	73,830	592	369,326	132
Shopping	800	193,377	225,129	750	1,185,625	1,022
Society	1,374	173,725	350,634	861	2,858,201	3,528
Sports	808	109,524	221,442	788	1,067,233	1,032
Technology	1,284	158,115	230,200	804	1,278,332	1,251
Travel	203	7,355	27,874	559	137,966	868
Weather	64	37,153	59,273	394	187,578	72

1.8.2 Propensity Score Matching Design

We employ the same propensity score matching technique as before, but now we create treatment groups based on session, task and domain instead of just task. We look across these three levels so that we can isolate the effect of each using Rubin’s causal framework and propensity score matching. This allows us to test if it is indeed task progression and not session or domain effects that is driving ad engagement. Since we stratify on session, task and domain combinations, we now have seven treatment groups to compare first (F), middle (M) and last (L) steps, i.e. {FFF, MFF, MMF, MMM, MML, MLL, LLL}. For example, MFF represents a middle step in the session but the first step of the task and the first step of the domain. We use the MMM observations as the control group as it is the middle step for each of the session, task and domain constructs in order to create 6 matched data sets for the test groups.

For each test group, we proceed by estimating the propensity score for each observation being in the treatment group relative to the control group. Then we again match each observation in the treatment group to an observation in the control group based on the linearized propensity score. This again creates a balanced set of observations for each treatment group. This gives us six balanced data sets for each session-task progression-domain combination.

Finally, we note that propensity score matching for such a large data set is a computationally expensive procedure requiring one to iterate through two lists of scored observations until at least one list is empty. We open-sourced¹³ a binary trees implementation of scored list matching which improved the computational performance from $O(N!)$ to $O(N\log(N))$ where N is the number of observations in the smallest treatment group.

1.8.3 Results

We again estimate a logit model with random effects for campaign ID, domain and user ID. We additionally control for ad position (above/below the fold or unknown), ad sizes, and device type (desktop/phone/tablet) as categorical variables. We omit day of week and time of day variables as they are again not statistically significant.

After fitting this model separately for each of the six treatment groups using the balanced datasets we can interpret the task step coefficient as the estimate of the average effect of task progression on ad click-through rates. Table 1-5 below shows the logit coefficient values for the treatment variables, the corresponding p-values and the odds ratios relative to the middle step by treatment group that these coefficient estimates imply.

¹³ <https://github.com/msdels/Matching>

Table 1-5 Data Set 2 Propensity Score Matching Model Output

Treatment Group (Session-Task-Domain)	Coefficient	p-value	Odds of Clicking	n	power
FFF	0.788	0	2.199	2,338,216	1
MFF	0.761	0	2.140	3,068,146	1
MMF	0.248	0.722	1.281	102,915	0.803
MML	0.286	0.096	1.331	363,489	1
MLL	0.107	0.001	1.113	2,643,893	1
LLL	0.453	0	1.573	1,699,535	1

Note: Bold numbers stand for significant at the 99% level.

Table 1-5 allows us to tease out the differential impact of task progression from that of pure browsing progression (session or domain) on ad engagement. We begin by looking at the middle session and first and last task step effects (MFF and MLL) as these most directly relate to our hypothesis. The logit coefficient for the first task step (MFF of 0.761) implies a 114% increased chance of clicking on a display ad¹⁴. Similarly, the logit coefficient of the last task step (MLL of 0.107) implies an 11.3% increased chance of clicking on a display ad. These results qualitatively replicate the findings in dataset 1 that first and last task steps lead to increased display ad engagement relative to middle steps. We also note that our estimates are sufficiently powered with all but the MMF case (with power of 0.8) reaching a power of almost 1.

The FFF coefficient represents the combined effects of first steps in the session, task, and domain. The first session step coefficient (FFF of 0.788) implies a 120% increased chance of clicking on a display ad, which is almost identical to MFF. To estimate the theoretical effect of starting a session (without starting a new task or domain) we take the difference of FFF and MFF. This effect size is not significantly different from 0, implying that it is the beginning of the task or

¹⁴ We note that these effects can be much larger for static display ad clicks (data set 2) than for video views (data set 1) but that might just be due to the average click-through rate being substantially lower than the average view-through rate (0.1% vs. 30%).

the domain and not the session that drives display ad clicks. Similarly, the last session step coefficient (LLL of 0.453) implies a 57% increased chance of clicking on a display ad relative to a middle step (MMM). By again differencing LLL and MLL we estimate the end of session effect as increasing display ad clicks by 41 percentage points above the end of task effect. These comparisons are all statistically significant at the 99% level. The first domain step (MMF) is not statistically significant, while the last domain step (MML) is marginally significant. Supporting results of a comparable model on the unbalanced data is shown in Appendix A.

Given that MML is marginally significant and both coefficients for MMF and MML are positive, we estimate a second model pooling first and last steps together for each of session, task, and domain partitions. This model increases the precision of our parameter estimates but it also assumes that both first and last step effects are of equal magnitude for each of the partitions. The results are shown in Table 1-6 below. Now the domain effect is more significant with a p-value of 0.056 and a coefficient implying a 36% increase in likelihood to click on an ad when transition between domains. The task effect remains significant with a coefficient implying a 69% increase in probability of clicking on an ad when transition between tasks. Given that task and domain effects are confounded, we can subtract the domain effect and conclude that the task effect alone implies a 33 percentage point increase in probability of clicking on an ad when transition between tasks above the domain effect. Following a similar logic, we can conclude that the session effect alone implies a 28 percentage point increase in probability of clicking on an ad when transition between tasks above the task and domain effects.

Table 1-6 Grouped Treatment Effects

Treatment Group	Coefficient	p-value	Odds of Clicking	n	power
Session Effect	0.673	0	1.960	4,037,751	1
Task Effect	0.527	0	1.694	5,712,039	1
Domain Effect	0.308	0.056	1.361	466,404	1

Note: Bold numbers stand for significant at the 99% level.

Collectively, these results support our hypothesis that task progression is indeed a significant determinant of display ad click rates after controlling for other relevant variables associated with ad engagement.

1.9 Validation Simulation for Advertisers

We now present a simulation to illustrate how it is possible to predict individuals' task progression in order to incorporate it into to display ad targeting. We assume that advertisers measure campaign success using a performance metric such as clicks or views and attempt to maximize it for each ad campaign using various strategies. They often operationalize this by constructing a predictive model over a set of features, usually some combination of user information, domain, time of day and page content. This modeling effort requires vast amounts of historical data to which a binary classification model such as logistic regression is calibrated (Perlich et al. 2012).

We take advantage of user navigation variables such as domain, page depth and page indicator to evaluate how incorporating the task progression construct compares to other standard targeting variables. We compare two logistic regression models using the variables available to us which allows us to evaluate the explanatory power of task progression for display ad targeting.

For this comparison, we return to data set 1 as we have full view of the domain visits and we again use view-throughs as our outcome variable. We partition the data by randomly selecting 50% of the users to the training set and the remaining to the test set. Using the training set we construct the same propensity scoring model as before to predict the probability of a given page being a first or last task step for each one of the users. Using this model, we create the task progression variables for the test set. They represent the estimated probability of a particular page view being a first or last task step for each person tracked. Next, we compare two linear models. First, we use the ad engagement model of the previous section with random effects as before to estimate the probability of fully viewing a video ad, and include our estimated first and last step task progression variables. Second, we estimate a standard logistic regression using only the first and last step task progression variables. The results are shown in Table 1-6 below.

Table 1-7 Advertiser Performance Comparison

	<i>Dependent variable:</i>	
	Video Completed	
	<i>random effects</i>	<i>logistic</i>
First Step	1.418***	2.730***
Last Step	0.464***	0.958***
Page Depth	-0.082	
Page Indicator	-0.375***	
Task step	-0.004	
Task step squared	-0.005	
www.economiamatin.fr	0.713***	
www.independent.co.uk	0.816***	
www.livefoot.fr	0.304	
www.remedes.de.grand.mere.com	0.632***	
Campaign Type: CPG	-0.629**	
Campaign Type: Electronics	-0.568**	
Campaign Type: Entertainment	-0.301*	
Campaign Type: Supermarket	-0.324*	
Campaign Type: Telecom	-0.195	
Campaign Type: Unknown	-0.374	
Campaign Type: Car	-0.305**	
Constant	-0.621**	-1.095***
Observations	9,696	9,696
Akaike Inf. Crit.	12,557.090	12,969.680

Note:

*p<0.1; **p<0.05; ***p<0.01

We first examine the full linear regression model with random effects. This model finds that both the first and last task step estimates are positive and statistically significant (at p<0.01)

and can be interpreted as predicting a 312% and 59% increase in the chance of completing a video relatively to predicted middle steps, respectively. Note that the coefficient estimates are also quite large.

Perhaps a more interesting comparison is the simple logistic regression using only the task progression variables with no random effects. No video view-through performance data is required for these predictions, and this level of performance can be achieved without having served a single ad. This model also manages to achieve a higher AIC than the more complex model (12,970 vs 12,557). Additionally, the out of sample AUC scores for the full model and just a linear combination of the task progression variables are both 0.619, indicating the model prediction is just as good as the full model. This is a significant step forward in addressing the cold start problem in predictive modeling for display advertising whereby practitioners have no historical data to rely upon in order to train their targeting models (Pan et al. 2019). Traditionally, this has been addressed through A/B testing and multi-armed bandit approaches, which require a substantial number of ads be spent on discovering high performing placements (Scott 2010). The number of ads needed for discovery increases linearly with the number of websites as well as the number of pages on the website and thus quickly becomes intractable. Display ad targeting using task progression provides a cost-free means of significantly improving baseline performance for cold start optimizations in display advertising.

To demonstrate the practical role that task progression can play in targeting, we return to the news task in the second data set since it is large enough to answer the question: ‘How many ads spaces need to be purchased in order for the click-based logit model to outperform the task progression model based on only observing web traffic?’ We apply the logit model as described earlier in this section and vary the number of ads spaces observed from 5,000 to 1,000,000 ad

spaces and then make out of sample predictions on equivalently sized sets of ad spaces. For task progression we only observed the opportunities for 5,000 randomly selected ad spaces and no performance data. We then repeat this process using 1,000 random samples.

To measure performance, we calculate the AUC score for the out-of-sample predictions adjusted for the cost of achieving such performance. The cost of achieving these performance numbers is the amount spent to purchase the ad spaces and the computing resources to store, process and predict using the click-based logit models. Since the task progression model only uses 5,000 random ad space opportunities (not bids, purchases, or clicks), it requires negligible computing resources and any modern laptop can easily perform these calculations for thousands of advertisers. However, click based data requires a large amount of computing resources. We estimate both costs as a percentage of revenue for three public DSPs at the time the data was collected. We then divided the observed AUCs by the average cost percentages of these companies (see Appendix A for details). The results are shown in Figure 1-5.

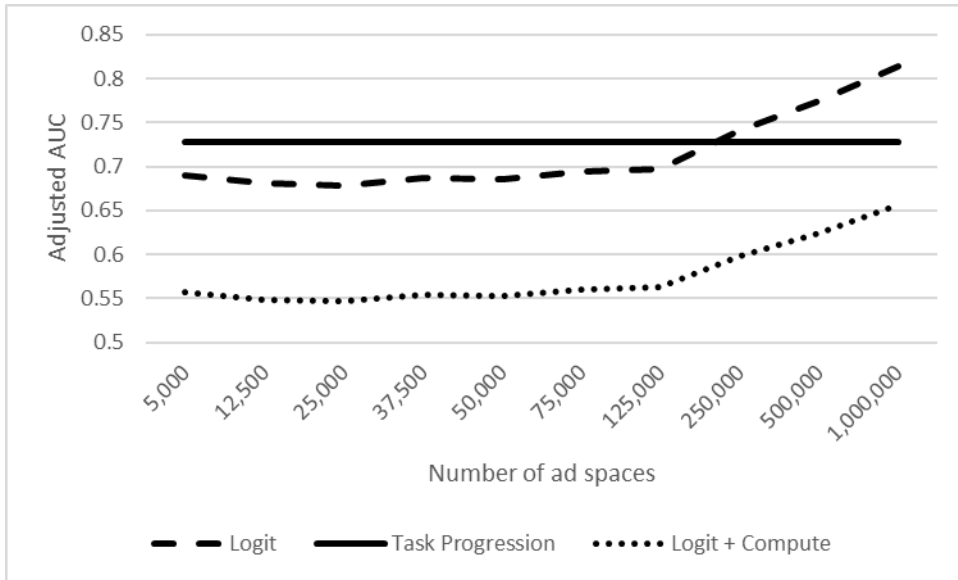


Figure 1-5 Price Adjusted AUC

From the simulation analysis we see that if we disregard the cost of computing resources, the click-based models will eventually outperform the task progression model after 250,000 purchased ad spaces. However, once accounting for the computing resources required to run and maintain click-based models, it becomes clear that on a cost-basis, the task progression model dominates the more data intensive click-based models.

Additionally, for the two weeks of observed RTB data, only 0.4% of domains had over 250,000 impressions. These domains accounted for 76% of ads bought and illustrate how skew the distribution of internet ad real estate is and shows how difficult it is to accurately score ad space in this long tail. This long tail is precisely why the cold-start problem in digital advertising (Pan et al 2019). Our simple model of task progression has demonstrated its ability to improve predictive accuracy on these less frequently trafficked domains. In particular, many domains do not get 250,000 total impressions in a month and a click-based logit model will never be able to score these ad spaces as well as a task progression model can.

Given that the average RTB CPMs range from \$0.50 to \$2 (Blaustein 2017), these 24% of total ads bought in the tail represent between \$123,000 and \$494,000 of ad purchase costs in two weeks. It is also possible that firms may increase their relative purchases of ad spaces in the long tail if they can better predict ad performance as this more accurate scoring may provide them with a competitive advantage.

1.10 Conclusion

Online advertising has come a long way since the popularization of the internet. Search advertising has worked incredibly well and has dominated much of digital advertising while display advertising, which most closely resembles traditional print media, struggles to perform well (Li and Kannan 2014).

We propose that a critical dimension overlooked in display advertising is knowledge of the user's task progression and their attentional capacity towards achieving their online browsing goals. In this paper, we explore consumers' web browsing tasks. Since many people browse the internet deliberately undertaking a task, we use website category information to classify these online tasks. Building on the marketing and attention literature, we hypothesize and show evidence that engagement with online ads is higher in the beginning and ending of tasks relative to the middle. Counter to previous display ad research showing either constant or diminishing attention over time, we provide evidence that ad engagement is highest at times when task engagement is lowest; i.e., at the task beginning and end points.

Using two datasets, we develop a novel causal inference application utilizing propensity score matching and random forests that exploit the hierarchical structure of websites. This provides us with a balanced data set to apply a standard random effects logit model to measure how task

progression influences the attention to ads. As hypothesized, we find that significantly higher levels of ad engagement occur at these task transition points. Our main result helps to explain that, much of display advertising is misaligned with the consumer’s attention to a primary task, leading it to being ignored.

1.10.1 Caveats

Given the nature of website structures needed for our inference procedure, we are not able to generalize our findings to unlimited feed-type sites like Twitter and Facebook. We do note that this is a measurement issue and that nothing in our theory suggests the same attentional process are not at play. A second important caveat is that we only examine measures of ad attention and not purchase behavior driven by display advertising. We encourage future research to build on this model by examining deeper funnel activities and other forms of digital advertising, although as mentioned earlier, this will require sophisticated attribution models. Finally, we do not investigate the potential interaction between task progression and behavioral targeting. Task progression predicts when users are likely to attend to an ad while behavioral targeting selects users for whom the ad is more relevant. Further research is needed to establish whether or not these two strategies are complements or substitutes.

1.10.2 Implications of Findings

Methodologically, our paper has implications for empirically modeling of clickstream data. After using propensity score matching to create a quasi-experimental design, we are able to show the parameter estimation bias arising from regression models estimated with standard observational datasets. We argue that page visits are not independent units of analysis since ad engagement is affected by task progression and that this should be corrected in any empirical modeling of such data. We also note that we provide a relatively unsophisticated model to predict

end of task and that it does not exploit any user level information or more complex website information (e.g. Trusov et al. 2016). Our modeling effort is sufficient for our purposes, but from a practitioner's standpoint, we provide a framework for targeting using task progression where real-world gains could be much higher than we estimated.

Substantively, our results also have important managerial implications for web publishers and digital advertising firms alike as they provide greater insight into the online attention process. Our results highlight when not to advertise, thereby allowing web publishers to reduce the ad load on their websites in the web pages where consumers are highly attentive to the online task at hand without losing much ad engagement. This will both increase the engagement rates to their other ads and make it a more pleasant experience for the user by reducing ad clutter. Additionally, Janiszewski et al. (2013) have shown that this type of ad reduction leads to a virtuous circle of increased ad efficacy over time. This, in turn, should lead to higher prices being paid for the publisher's high attention-drawing ad spaces. In regard to advertisers, they should incorporate task progression as an input into their ad buying and placement strategies, particular for less trafficked domains where direct performance data is scarce. In a world of ever increasing big data, advertisers have a growing amount of opportunity to learn and predict task progression. By ignoring the task, advertisers are ignoring the attention capacity of their targeted consumers and are potentially wasting a large amount of ad dollars.

In sum, this paper is a first attempt at showing how task progression affects attention to and engagement with display advertising. By incorporating this construct into online display advertising strategies, advertisers can increase the effectiveness of this important marketing tool. If this is done in a rigorous manner, we believe all interested parties, advertisers, publishers and consumers will collectively benefit.

2 Real-time Digital Ad Allocation: A Fair Streaming Allocation

Mechanism

2.1 Introduction

Digital advertising spending continues to grow as a share of the marketing mix and is predicted to reach \$129B in 2019, it is also predicted to overtake traditional ad spending (Ha 2019). While firms like Google and Facebook enjoy large portions of this growth, new digital channels like programmatic advertising have flourished as well. While much research has focused on search and social media, little marketing research has studied the new challenges in programmatic advertising.

Programmatic advertising is projected to make up almost half (\$57B) of all digital advertising spending in 2019 (Fisher 2018). Programmatic advertising is defined as the use of automation in the buying, selling, or fulfillment of digital advertising. Programmatic advertising takes place on real-time bidding (RTB) exchanges where publishers list ad spaces in real-time as users load webpages and then advertisers can bid on these ad spaces. Due to the technical barriers to entry of the marketplace, advertisers have outsourced the procurement of ad space to demand-side platforms (DSPs) who bid on their behalf.

In previous research, Balseiro et al. (2014) investigate the supply side of the programmatic market. From a publisher perspective, the problem is one of profit maximization where they must choose whether or not to list their ad space on an RTB exchange or not and at what price. These authors provide an optimization strategy where the ad spaces are always listed on the RTB exchange with a dynamic reservation price. The reservation price is adjusted based on the

publisher's projected ability to meet their own direct sales contracts. They show how using simple greedy allocation fails to allocate ad space efficiently and can lead to losses in yield of up to 70%.

The demand side of the programmatic market is substantially different from the previously studied adwords and publisher problems. Here, agencies and advertisers contract with DSPs to procure a set number of ad spaces on their behalf. DSPs, therefore, face a steady stream of ad spaces being supplied by the RTB exchanges. At this point, they must in real-time: (1) score these ad spaces, (2) decide which of those ad spaces they wish to purchase, (3) which advertiser to allocate it to, and (4) how much to spend on the ad space. Additionally, this supply of ad space is highly variable making it infeasible to rely on traditional models with stationary distributions that attempt to learn optimal strategies based on a training set as common in most similar problems

In this paper, we investigate the demand side of the programmatic market. We introduce the DSP problem and differentiate it from the closely related publisher and adwords problems (both well described in Mehta 2012). In brief, the adwords problem is also a real-time allocation problem, in the case of adwords for a search engine that receives keywords and bid amounts from advertisers along with a monetary budget. The search engine wishes to allocate the ads spaces to maximize the sum of budgets spent across all advertisers. Whereas, the publisher problem refers to website publishers that presell their ad space to various advertisers based on expected future user traffic. Their optimization problem is to allocate incoming ad spaces to their various advertisers in such a way that it fills their obligated quotas and maximizes a pre-set performance metric (e.g., clicks).

We show how a DSP should decide, in real-time, which ad spaces to bid on and which advertisers to allocate them to. Similar to Balseiro et al. (2014), we show that careful capacity management is crucial and provide an online algorithm that performs significantly better than

greedy allocation and alternative methods, including benchmark methods for the adwords and publisher problems. Through both real-world exchange data and simulated data, we demonstrate that our algorithm is superior to comparable algorithms and that it can perform similarly to offline methods such as MacAfee’s (1994) Alternating Selection Mechanism.

This paper is organized as follows. We first conduct a brief literature review from both computer science and economics of matching mechanisms in our context. We then describe our proposed algorithm before moving onto evaluating our method against several competing methods. In the final section, we conclude.

2.2 Literature Review

The research problem we examine in this paper is a variant of other well-studied online ad allocation problems. It is related to both the adwords problem and the publisher problem. The adwords problem is a search engine revenue maximization problem. In this setting, the advertisers approach the search engine and provide a daily budget along with a bid price for individual keywords. The search engine’s goal is to exhaust as much of the budgets as possible. Most current solutions assume that searches arrive in random order but tend not to be too dissimilar from the previous days allowing reasonably effective strategies to be learned from historical data.

More formally, Kalyanasundaram and Pruhs (2000) define the adwords problem through a bipartite graph with one set of vertices U to represent advertisers and one set of vertices V to represent user searches. Each advertiser $u \in U$ has a budget $B_u > 0$, and there is a set of edges $(u,v) \in E$, with an edge to represent a bid by u on search v and annotated with weight $bid_{uv} > 0$, denoting the bid amount. When a search query $v \in V$ arrives, it needs to be matched to some neighboring advertiser $u \in U$ who has not yet spent all its budget. Once we match v to u , then u depletes bid_{uv}

amount of budget. If a vertex u finishes its entire budget, it becomes unavailable. The goal is to maximize the total money spent.

In the publisher (or display advertiser) problem, the publisher has quota-based contracts from multiple advertisers for a set number of ads to place every day. The goal of the publisher is to primarily satisfy these quotas and then secondly to perform adequately on a key performance indicator (KPI) such as clicks; therefore creating more of a covering problem instead of a packing problem. Users visit a publisher's site in random order, and therefore, the ad spaces appear in random order as well. These solutions also tend to assume that publisher site visits are unlikely to vary greatly from day-to-day, allowing reasonably effective strategies to be learned from historical data.

Dimitrov and Plaxton (2008) more formally defines the publisher problem through a bipartite graph where one set of vertices U represent advertisers and one set of vertices V represents different publisher sites. The edges of the graph $(u,v) \in E$ have weights w_{uv} , representing a KPI score, and the vertices $u \in U$ have capacities c_u . As before, when a vertex in V arrives, representing a user visit to a publisher site, it has to be matched to a neighboring advertiser in U , such that each $u \in U$ is matched at most c_u times. The goal is to maximize the total weight of the matched edges.

Both problems are discussed at length in the computer science literature (e.g., Devanur and Hayes 2009, Manshadi et al. 2010, Gollapudi and Panigrahi 2014, Feldman et al. 2018). These algorithms all represent variations of algorithms for online bipartite graph matching. The proposed solutions tend to set up and solve the problem using primal-dual linear programming. Most treatments of this problem focus on short-term revenue maximization only. Any degree of fairness or equality is achieved either through some minimum amount of random assignment or minimum coverage targets for quota fulfillment (Gollapudi and Panigrahi 2014). The methods require an

offline training period where the assumption is that the distribution of arriving items (keywords or ad spaces) remains reasonably stable after that. A more dynamic environment is handled by updating the training set more often and repeating the learning procedure. As mentioned earlier, the DSP version of this problem is likely to be more variable, and assuming such stability in the solution concept is not practical. To evaluate our proposed method, we shall compare our results to the Devanur and Hayes (2009) and Feldman et al. (2018) primal-dual formulations of the adwords and publisher problems, respectively.

The economics literature has also provided much insight into this class of problems, and in particular, has provided reasonable offline solutions. It is worth noting that much of the economics literature is far more concerned about eliciting truthful responses from advertisers; however, they also provide valuable design principles and pricing frameworks.

Early work by McAfee (1992) describes three auction designs to solve this problem for two competing agents where the auctioneer does not know the values of the items. The first design is the Winner's Bid Auction (WBA) where two agents are each shown an item, they declare their bid in private, and the highest bidder receives the item and pays half of their bid amount to the loser. The second design is the cake-cutting mechanism (CCM). In this mechanism, one agent proposes a price, and the other agent either takes the good or accepts the price. The third design is the Alternating Selection Mechanism (ASM). In this design, all the items are presented at once, and the two agents alternate in selecting their most valued item until all items are gone. The ASM is shown to be remarkably efficient, where under general assumptions, it will have an efficiency loss that is at most $\frac{1}{2n}$ times the highest value item, where n is the number of agents.

While these three mechanisms do not directly address any kind of budget or constraint, they do assume that each agent has a right to each item and they provide mechanisms for fairly allocating these items. Of note is the ASM mechanism, which, unlike CCM and WBA, does not rely on monetary transfers and is solving a very similar problem to the second part of our proposed solution. The major drawback with ASM is that it is not implementable in an online fashion and therefore, cannot directly solve our problem. Despite this drawback, these methods of distributing items fairly are insightful. The critical intuition being that fairer allocations can be achieved by agents alternating when they get items and by some form of value transfers to compensate agents for not getting an item. Additionally, ASM provides a useful benchmark for us in that it is a well-studied solution to the offline version of our problem.

Whereas McAfee (1992) looks at allocation rules for a known set of items, Casella (2003) investigates voting mechanisms where agents also face an uncertain number of future items that they may care about differently. They show that for agents facing an uncertain stream of votes, it can be optimal for the agent to forgo voting on a present low valued item in order to be bid more on a higher valued future item. Additionally, they show that in general, welfare gains hold for markets with more than two agents. This mechanism is useful for voting schemes where all agents share the items, for example in a setting with a public good, but it fails to generalize to the case of indivisible private goods. However, the intuition of using votes as a scrip currency and passing on less desirable items to bid more on more desirable futures is useful in designing mechanisms for the present context.

Instead of looking at an unknown stream of future items, Jackson and Sonnenschein (2005) look at a stream of items from a known distribution. They propose a linking mechanism for a Bayesian collective decision problem where preferences of agents are private. By applying the law

of large numbers, they demonstrate that by linking many equivalent problems over time that the mechanism can elicit honest preference reports from the agents. This mechanism can be applied to many general social choice functions where the distribution of items is known and effectively solved the incentive problem without the need for monetary transfers. However, the problem faced by the DSP is the inverse of this, in our setting, there are no incentive problems since the preferences are always known, but the distribution of items is not.

2.3 Fairness

We briefly digress in this section to make an argument for fair allocations within the DSP. While fairness is often a consideration in the computer science literature, it is much less discussed in the business literature. Traditionally, business researchers are interested in either profit or revenue maximization. However, given the rising distrust in technology companies and greater calls for transparency in algorithms (e.g., "Algorithm and blues," 2016), we argue that there is a strong case for fair algorithm design in matching markets and that firms' long-term sustainability could depend on algorithmic fairness. It is no longer uncommon for executives from large technology companies to testify before Congress (e.g., Facebook and Google) regarding algorithmic transparency, and it now seems reasonable to assume that many algorithms will one day be publicly exposed. It should, therefore, be a strategic consideration for companies to design fairer algorithms that would not cause a large backlash from end-consumers or business partners should they one day be made public. Fairness can also be seen as a good business practice that will keep their various partners happy and build trust, which in turn, will ensure a long term partnership.

For our purposes, when we discuss "fair," we mean that, relative to each other, advertisers are treated fairly within the DSP. However, fairness is an intrinsically relative measure, unlike efficiency, which is an additive measure. The DSP faces many ad spaces daily that are each desired

by multiple advertisers and they need to choose some mechanism to allocate these ad spaces. The problem the DSP faces is that it needs to choose between advertisers immediately and without accurate future knowledge of ad spaces. A greedy assignment would maximize only a coarse measure of short-term efficiency. We say coarse because the scores themselves are noisy and hard to compare across advertisers. We also say they are short term because, as mentioned previously, consistently poor performance for some advertisers will result in less future business and therefore cannot lead to long-term efficiency.

The opposite of a greedy assignment is a random assignment. This is synonymous with claiming that all advertisers have an equal right to an ad space and that each stands an equal chance of receiving a contested ad space. While certainly fair, it is far from efficient. It is not difficult to see that on many occasions one advertiser scoring the ad poorly will randomly receive that ad space at the expense of another advertiser scoring it highly.

Fair allocation in these online allocation problems has remained somewhat elusive. Devanur and Hayes (2009) acknowledge its importance and difficulty but offer a greedy algorithm while deferring fairness to future research. Gollapudi and Panigrahi (2014) propose a Max-min fairness algorithm that seeks to first satisfy a minimum baseline for all advertisers before turning greedy. While a step in the right direction, this class of methods still imposes binding constraints on hard to satisfy advertisers that can come at a considerable cost to overall efficiency.

We propose that fairness should only be considered for advertisers who demand an item and only on an item-by-item basis. Rather than assigning the item to the advertiser with the highest value, we borrow from the linking literature and propose assigning it to the advertiser in proportion to their demand for it relative to other advertisers' demand. Unlike in linking where each agent gets an integer count for the number of times they can bid a high value, we propose that we assign

items using a form of fractional assignment. This will lead to advertisers getting the item a proportional fair amount of the time in expectation. Additionally, linking this expectation to a scrips budget pushes towards fairer outcomes faster than pure random assignment would.

2.4 Algorithm Design

2.4.1 Environment

Before describing the algorithm, we first discuss the real-time bidding (RTB) advertising environment in more detail. This is a marketplace where publishers list ad spaces that are bid on by various demand-side platforms (DSPs) who represent advertiser interests. Ad space is listed in real-time, that is, as a user loads a publisher's webpage, the publisher sends a request to the RTB exchange to list the ad space. These requests are forwarded to the DSPs who are given up to 100ms to bid on this ad space using a blind second-price auction with reserve prices. Given the speed and resources required to operate in this market, advertisers contract DSPs to purchase ad space on their behalf algorithmically.

While advertiser contracts with DSPs primarily focus on ad impression quotas, they also typically include a key performance measure (KPI) that the advertisers are trying to maximize (e.g., number of clicks or purchases) over a set number of ads (usually a daily ad impression quota) (Zhang et al. 2016). The goal of the DSP is to keep all of its clients, that is, first meet the impression quota and then the performance standards. DSPs run a variety of proprietary algorithms that score all incoming ad spaces for each advertiser and their prespecified KPI. Since DSPs have many contracted advertisers (often thousands), many of the incoming bid requests are desired by multiple advertisers (i.e., the various ad campaigns). The DSP can only assign the ad to a single ad

campaign. Therefore, there is a tension between assigning an ad to the campaign scoring it highest (i.e., greedy allocation) versus a more equitable distribution of ad space.

The DSP problem is similar to the publisher problem, but with two unusual characteristics: (1) the distribution of ad scores and ad volume is highly unpredictable, and (2) there are no incentive problems for the advertisers (i.e., the advertiser ad campaigns are internal to the DSP) to report untruthful preferences for the ad spaces as these scores are provided by the DSP.

In general, we cannot learn the distribution of the arriving ad space scores for at least four reasons. First, the exchanges themselves route ad requests in opaque ways due to internal load balancing. Second, significant fluctuations in scores arise through regular operating issues such as server outages, internet connection problems, and software updates. Third, DSPs update their scoring algorithms regularly, often separately for each campaign, and sometimes in real-time for behavioral targeting. Finally, publishers themselves often sell ad space directly and can suddenly change their supply or quality of ads sold on the exchange. Jointly, these factors make the matching problem faced by the DSP less predictable than the well-studied publisher-side advertising problem or the adwords problem.

The second unusual market characteristic is that a DSP assigns all the campaign-ad scores based on proprietary models. Unlike most other matching problems, these scores are not private information of the advertiser but calculated by the DSP itself. It is, therefore, a given that the reported scores to any assignment mechanism will always be truthful.

At first glance, it may seem optimal to assign the incoming ad spaces to the campaigns scoring them the highest (i.e., greedy assignment). This is suboptimal for the DSP two reasons. First, scoring models are KPI- and campaign- specific, and since internal scoring models can be

very different, some campaigns tend to have higher scores than other campaigns. This can lead to near stochastic dominance in score distributions of some campaigns over others. This in itself is not a problem, but the scores for each ad space are correlated across campaigns. Thus, a greedy assignment mechanism would allocate many of the top-scoring ads to the top subset of campaigns. Similarly, a bottom subset of campaigns will receive almost entirely low scoring ad spaces leading to the eventual loss of business. This could result in some advertisers experiencing extended periods of sub-par campaign performance and ultimately the cancellation of their DSP contracts. Therefore, the optimal short-run allocation of a greedy algorithm does not lead to the optimal long-run allocation for a DSP.

A second concern is that these scoring models are quite narrow in scope. They often only concern the direct KPI and historical performance. Advertisers ultimately care about revenue but use a click KPI as a proxy for ad engagement and future purchases. Focusing on a single ad space with the highest click rates is not necessarily going to lead to the most revenue, especially since this creates substantial publisher incentives to game the KPIs or engage in ad fraud. Furthermore, the scores are noisy estimates of performance. By increasing the variety of ad spaces purchased, the DSP reduces the variance on expected performance, increases their market reach, and reduces the risk of ad wearout. Additionally, ad diversification lowers the risk of ad fraud.

Alternatively, one could assign the contested ad space randomly. This would be closely related to a serial dictator assignment (Abdulkadiroglu and Sonmez 1998) and can be seen as most 'equitable.' However, it is not difficult to imagine circumstances where an ad space goes to a campaign with a very low score, and another campaign with a very high score misses it. The mechanism proposed in this paper allows a sliding scale between random assignment and greedy assignment by allowing a weighting based on relative preferences.

2.4.2 Proposed Solution

In the next section, we proceed with defining an online allocation mechanism to solve the DSP problem. The mechanism has two broad components, an agent level thresholding algorithm and a market level allocation mechanism. The DSPs need a fast decision rule for deciding which items to compete for on behalf of an advertiser (or agent). This is the role of the agent-level threshold rule. The market level allocation mechanism is a means of allocating an item to one of the competing advertisers in the event that multiple advertisers desire the same item. After a closer look, this is not dissimilar from how the previous literature that treats these two problems through a single algorithm. For example, in the primal-dual formulations (e.g., Devanur and Hayes 2009, Gollapudi and Panigrahi 2014, and Feldman et al. 2018) we effectively also have a thresholding rule where the observed ad space score is adjusted by a learned weight and only gets bid on if it is above 0. Secondly, these methods assign the ad space to the largest adjusted score above 0. This adjustment is based on fairness or capacity criteria, making it different from greedy but also comparable to our market-level optimization. Our method handles fairness explicitly in the market level step by tracking items won and lost over time.

2.4.3 Part 1 - Agent level optimization

We begin with the agent level optimization, and to gain intuition into the solution, we focus on a DSP with only a single advertiser (agent). The DSPs are tasked with smoothing over the ad supply variability by procuring a fixed number of ads for their clients each day, typically spread out across the day in some pre-specified manner. The DSP is incentivized to buy the best ads possible for their clients subject to meeting their daily quota requirements. However, since ad space availability on the exchanges is so variable, it is hard to come up with a static set of buying rules once-off that will both meet the quota and allocate the best ads possible. The problem for the

DSP is that they do not know what ads will become available throughout the day and how they should be buying optimally throughout the day for even a single advertiser.

To illustrate the problem, consider Table 2-1 below. Here we see the sequence of ad opportunities that appear one-at-a-time with their corresponding scores. A decision must be made on Ad1 before Ad2 is revealed. Suppose that this sequence represents three time intervals of $\{(Ad1, Ad2, Ad3), (Ad4, Ad5, Ad6), (Ad7, Ad8, Ad9)\}$ and that the DSP would like to purchase one ad in each time period. With perfect hindsight, setting a threshold of $T=7$ would succeed in buying only the single highest scoring ad in each time period.

Table 2-1 Single Advertiser Ad Sequence

	Ad1	Ad2	Ad3	Ad4	Ad5	Ad6	Ad7	Ad8	Ad9
Score	7	3	2	4	8	1	3	9	6

More formally, we can define the DSP problem through a bipartite graph where one set of vertices U represent a single advertiser and one set of vertices V represents different ad spaces appearing on an RTB exchange. The edges of the graph $(u,v) \in E$ have weights w_{uv} , representing a KPI score, and the vertices $u_t \in U$ have time period capacities c_t . When a vertex in V arrives, representing a user visit to a publisher site, in time period t it can be ignored or matched to a neighboring capacity slot in u_t , such that each u_t is matched c_t times. The goal is to maximize the total weight of the matched edges.

To make this problem more tractable, we assume that bid prices are external to this mechanism and that the buying firm can in real-time provide a scored preference each ad spaces for an advertiser. Additionally, we disregard the fact that not all bids are won, we focus only on the ad spots won on the exchange. This is not problematic as our proposed solution moves a

threshold in order to win enough ad space and does not directly depend on the auction win rate. If we win too many ads, the threshold needs to be raised. Conversely, if we win too few ads, then the threshold needs to be lowered. It also seems plausible that the buying firms can score ad slots in some way as this is precisely the nature of their business and well described in Perlich et al. (2012). Finally, we also assume that we are given how the advertiser would like their impressions to be spread across the day. It is often the case in practice that the daily distribution is provided by the advertiser directly or by another optimization algorithm focused on optimal ad timing.

In this setting, we wish to have a simple rule of thresholding where the advertiser accepts any item with a score above the time-dependent threshold $T(t)$ and rejects any item below it. Thresholding is a standard strategy in feedback control loops (Sung 2009), and it also provides a simple rule for the DSP to apply in real-time since the ad space auctions happen in under 100ms. What remains is to find an efficient way to adjust this threshold value to fulfill advertiser demand and adjust to any shocks in the system.

III.2 Smoothing

The first problem we encounter in this system is that the supply of ads is highly variable, making it difficult to estimate the number of ads available at even a static threshold. To overcome this, one cannot target the actual number of ads won but rather a smoothed estimate of the number of ad spaces we would expect to win at the current threshold. For this purpose, we propose using a Dynamic Linear Model (DLM) (West and Harrison 1997).

DLMs are fully Bayesian, and if conjugate priors are chosen, then the model is also completely online. Additionally, these models can also self-calibrate the noise component, and

learn seasonality and higher-order rates of change. For this paper, we use a standard DLM that reduces to the Kalman Filter (Kalman 1960).

$$x_k = F_k x_{k-1} + w_k$$

$$z_k = H_k x_k + v_k$$

Here x_k is the observed number of ads won in a time period k , w_k is a gaussian error with a learned variance term, and F_k is the transition matrix. z_k is the true (smoothed) ad rate for period k , v_k is also a Gaussian error with a learned variance term, and H_k is the observation matrix. F_k captures the evolution of the state space from one period to the next. In our case we naïvely expect x_k to remain unchanged and set F_k to the identity, but one could also learn cyclical trends and have F_k be time-dependent. Similarly, H_k captures how the hidden state realizes in the measurement z_k . In our case we again expect H_k to be the identity.

III.3 Optimal Control

Now that we have a smoothed ad rate to control, we need a control algorithm to adjust the threshold to get to the target ad rate that we desire. A standard tool from control theory is a proportional integral derivative (PID) controller Sung (2009). This allows us to adjust the threshold $T(t)$ based on the difference between the ad rate we are achieving and the desired ad rate. The setpoint (SP) is the desired number of ad spaces an advertiser wishes to win, and the process value is the observed number of ad spaces won. A proportional-integral-derivative controller is a control loop feedback mechanism that uses only the error term (PV-SP) of each time period to adjust the threshold. It uses the following parameterization:

$$T(t) = K_p e(t) + K_i \int_0^t e(\tau) d\tau + K_d \frac{d}{dt} e(t)$$

Table 2-2 PID Parameters

Parameter	Description
T(t)	Score threshold
SP	SetPoint (target ad rate)
PV	Process Value (smoothed ad rate)
E(t)	Error (PV-SP)
K _p	Proportional Gain
K _i	Integral Gain
K _d	Derivative Gain
t	time period
τ	Variable of integration; takes on values from time 0 to the present time period

The PID algorithm works in three parts, in the first component, K_p weights the threshold adjustment based proportionally on the latest error estimate. This is a first-order adjustment based on the most recent error information. Second is the integral component, it keeps track of all previous errors and weights the adjustment based on K_i times the accumulated error. This allows the adjustment to speed up if it sees a growing total error over time despite the incremental changes. The final term is the derivative component, K_d, which adjusts the threshold based on the change in the last two errors terms. The derivative component allows basic linear extrapolation of how effective the previous threshold adjustment was at reducing the error, and can be effective at dampening any potential overshoot in threshold adjustments.

III.4 Self-calibration

Finally, there are two general requirements for a PID controller to function correctly. The first is a target for the controller (SP). In this case it is the number of ads desired during a time unit interval. We can compare this with the number of ads purchased during the previous time units and decide whether or not the threshold needs to be increased or decreased. The second requirement is that PID controllers need to be tuned to the system in which they operate. This second requirement is non-trivial as these controllers can become highly unstable if tuned incorrectly for the system they are controlling.

For the first requirement, adjusting the threshold to hit a target is very difficult when the process value is so noisy. Even if we left T static for the entire day, each time unit would lead to a different number of ads being bid on because the number of ads observed with a score above T is random. This is precisely why we use a DLM to replace the observed ads with the prediction estimate for the state of the underlying ad rate. The DLM is also a model that can tune itself according to the noise in the system to provide an accurate estimate of the current state as well as an accurate prediction for the next time period.

For the second requirement, we need a learning strategy to tune the PID controller dynamically. PID controllers are very sensitive to calibration and require specific tuning for each application (Kim et al. 2008). For the proposed algorithm, we start the controller in a relatively sluggish and stable state. Then using particle swarm optimization (PSO) to test out new tuning parameters based on whether the size of the previous error was too big, too small or adequate (Kennedy and Eberhart 1995).

To do this, we use the Ziegler-Nichols (Ziegler and Nichols 1993) transform (see Table 2-2) of the PID parameters to change the parameter search space from 3 to 2 dimensions¹⁵. Thus, the parameter tuning now occurs over only K_u and P_u instead of K_p , K_i , and K_d .

Table 2-3 Ziegler-Nichols Transformation

	K_p	K_i	K_d
PID	$0.6K_u$	$2K_p/P_u$	$K_pP_u/8$

For PSO, at each time step, we create several new particles that represent slight variations of K_u and P_u . We perturb one particle at a time by a factor of $(1+\sigma)$ to create four sets of new candidate parameters, where σ is the step size of the perturbation. For each particle, we look at what threshold adjustment would have been in the current time step. We then calculate the proportional error of the current time step as $\epsilon_{prop(t)} = \frac{SP-PV}{SP}$. If $\epsilon_{prop(t)} < \delta$ then leave the PID parameters unchanged. Otherwise, we change the parameters to particles that would move the threshold most in the direction that error indicates. δ is the error tolerance we deem acceptable for not changing the tuning parameters. Both σ and δ should be learned during an upfront calibration period. This algorithm is illustrated below in Figure 2-1.

¹⁵ The PSO needs to search across parameter space to find the optimal calibration for the three PID tuning parameters. The Ziegler-Nichols transform provides a generally stable link between the three parameters that reduces them to two parameters. This simplifies the search process and decreases the compute time needed to execute the algorithm.

```

Initialize starting turning parameters  $K_{u,0}$  and  $P_{u,0}$ 
Set search tolerance  $\delta > 0$  and search step size  $\sigma > 0$ 
3: for each time step  $t$  do
    calculate the proportional error of current time step as  $\epsilon_{prop}(t) = \frac{SP-PV}{SP}$ 
    if  $|\epsilon_{prop}(t)| < \delta$  then
6:         leave  $K_{u,t+1}$  and  $P_{u,t+1}$  unchanged
    if  $|\epsilon_{prop}(t)| \geq \delta$  then
        Based on  $\epsilon_{prop}(t)$  evaluate if  $T(t)$  needs to be increased or decreased.
9:         Test alternate parameters in previous PID state:
             $K_{u,t+1}^1 = K_{u,t+1}(1 + \sigma)$  and  $P_{u,t+1}$ 
             $K_{u,t+1}^2 = K_{u,t+1}(1 - \sigma)$  and  $P_{u,t+1}$ 
12:          $K_{u,t+1}^3$  and  $P_{u,t+1} = P_{u,t+1}(1 + \sigma)$ 
             $K_{u,t+1}^4$  and  $P_{u,t+1} = P_{u,t+1}(1 - \sigma)$ 
        Select new  $K_{u,t+1}$  and  $P_{u,t+1}$  that move  $T(t)$  the most in the desired direction

```

Figure 2-1 Particle Swarm Implementation

This type of search method has the benefit of leaving the PID controller in a stable state once the system has found a parameter set well suited to the current environment. Conversely, during a system shock, a static PID may take a long time to adjust, so the PSO will re-calibrate the PID to become far more responsive. Then once the PID finds its new steady state, the PSO will again re-calibrate the PID to become less responsive.

2.4.4 Market level allocation

The agent level solution works well provided that no other advertiser wishes to buy to the same ad space. However, many ad spaces will appear where multiple advertisers score it above their agent level thresholds. It is only in these circumstances where an ad space is considered

contested that the market level allocation algorithm needs to be activated to fairly allocate the ad space to one of the competing advertisers.

To illustrate this problem, consider Table 2-4 below. As in Table 2-1, we observe the sequential scores but this time for multiple advertisers. In this scenario, Ad1 has passed the threshold for each of the three advertisers and needs to be assigned to one of them before Ad2 is revealed. Since the threshold has already accounted for undesired ad opportunities, all three advertisers desire the ad space. Therefore, the only concern at the market level is how to pick which of the competing advertisers should be assigned the ad opportunity.

Table 2-4 Multiple Advertiser Allocations

	Ad1	Ad2	Ad3	Ad4	Ad5	Ad6	Ad7	Ad8	Ad9
BMW	9	9	4	8	10	9	6	3	10
Mercedes	8	8	5	6	7	8	4	5	9
P&G	4	3	8	5	2	3	5	9	4

In this simple example we can see why a few standard solutions may not work. First, consider a greedy assignment mechanism where that ad opportunity is always assigned. Since BMW allows dominates Mercedes in score, all the contested ad opportunities will be assigned to BMW. This could lead to starving Mercedes for quality ad spaces. Additionally, any type of sequentially alternating mechanism would be no better than random assignment and lead to obvious inefficiencies in allocations. What is desired is a more strategic alternating scheme when inefficiencies are small and that only occur if an advertiser has been starved for a while. This is the precisely the intuition behind the proposed solution.

Next, we proceed to describe the allocation algorithm illustrated in Figure 2-2. Each advertiser is provided a score for all items $r_{ij} \in [0, 1]$ which is provided truthfully at the moment

the ad space j arrives. Each new advertiser i begins with an initial scrip budget of $B_i = 1$. The market level allocation is only relevant to advertisers that see an ad spot with a score above their PID threshold. We use A to denote the set of all the advertisers that would like to bid on item j . At each iteration, a new item j appears. We let b_i denote the bid that advertiser i is making for item j . b_i is only set for advertisers in A (i.e., advertisers that have r_{ij} above their PID threshold value), for the remaining advertisers it is equivalent to 0 since they abstain.

b_i is defined as a function of how desirable it is (r_{ij}) and how much budget they have (B_i). This is useful in two ways, first, the advertisers that value the item more spend more on it. Second, as we shall see shortly, B_i acts a fairness tracker, the higher its value, the more an advertiser bids on an item. The bidding function $f(r_{ij})$ illustrated in Figure 2-3 has a tuning parameter γ where $\gamma=1$ leaves the scores unchanged (i.e., it will return r_{ij}). The function leaves the highest score of all the advertisers unchanged, but then either moves the other advertisers closer to highest score ($\gamma<1$) or it moves the other advertisers' scores closer to 0 ($\gamma>1$). As gamma approaches 0, the algorithm will assign items randomly as all the scores will become identical. As gamma approaches infinity, the algorithm will assign items almost greedily since all but the top score will be pushed to 0. This allows for an easy sliding scale for DSPs to use to move between random and greedy assignment. We leave γ as a firm-specific tuning parameter to make the fairness/efficiency tradeoff and demonstrate in simulations in the next section how it can be used.

These bids (b_i values) are now used as weights to generate categorical probabilities for each advertiser's chance of being assigned item j . By assigning the item probabilistically, we further address the fairness concern of the advertiser with the highest rating always getting the item. In a single item setting, each advertiser will get the item in the relative proportion that they prefer it.

This becomes the discrete and online analog to proportionally fair allocation described by Cole et al. (2012), and similar to the ideal fractional allocation from Feldman et al. (2018).

We then draw the winning advertiser from this categorical distribution and allocate item j to them. This kind of probabilistic assignment, spread across multiple items, can be seen as practical means of achieving fractional assignment.

The winner's bid (b_{win}), is subtracted from their balance B_i . b_{win} is then redistributed between all the advertisers (including the winning advertiser) according to the proportions that they scored the item. This keeps the scrip supply constant and increases the losing advertisers' chances of winning an item later on. The increase in each advertiser's scrip is proportional to how much they value it relative to other advertisers. This storing of B_i can be seen as similar to a continuous version of Casella's (2003) vote storing idea. The redistribution of scrip is also very similar to Cramton et al.'s (1987) trading mechanism when the advertisers 'own' equal shares of an item but have different values. The main difference here is that the item is not always allocated to the advertiser with the highest value.

Lastly, to keep the scrip supply constant per advertiser, we re-normalize B when an advertiser leaves. Any new advertiser i entering starts with $B_i = 1$.

This mechanism can also represent ordinal preferences. By adjusting the advertisers' ordered preferences using a probability integral transform, as in Santos et al. (2016), we can get uniform item preference values to use in the above mechanism.

-
- 1: Initialize scrip budget: $\mathbf{B} = 1$
 - 2: Draw next item j
 - 3: Agents declare their preferences $r_{..} \in [0, 1]$
 - 4: **for** $i \in A$ **do**
 - 5: $b_i = f(r_{ij}^\gamma) \times B_i$
 - 6: Generate categorical probabilities: $p_i = \frac{b_i}{\sum_k b_k}$
 - 7: Draw winner i from $Cat(p_1, \dots, p_n)$ and allocate item
 - 8: Update $B_{win} = B_{win} - b_{win}$
 - 9: **for** $k \in A$ **do**
 - 10: Update $B_k = B_k + \frac{r_{kj}}{\sum_l r_{lj}} \times b_{win}$
 - 11: **if** Some agent l leaves the market **then**
 - 12: re-normalize \mathbf{B} : $\mathbf{B} = (N - 1) \frac{\mathbf{B}_{-l}}{\sum_{k \neq l} B_k}$

Figure 2-2 Market Level Algorithm

-
- Given preference matrix $R_{A \times I}$ and γ
- 2: **for** $i \in I$ **do**
 - for** $j \in A$ **do**
 - 4: $r'_{ij} = r_{ij}^\gamma$
 - $r'_i = r'_{i.} / \max(r'_{i.})$
 - 6: $r'_i = r'_i \times \max(r_{i.})$

Figure 2-3 Score Normalization

2.5 Evaluation

We now turn to evaluating the proposed algorithm. To our knowledge, no other algorithms describe the DSP problem that we present in this paper. Therefore, we analyze the proposed solutions in two steps to highlight each solution's part and show how they can be best compared to the best strategies available to DSPs. First, we begin by examining the agent level solution and

compare it with similar methods used in practice. We do this using a real-world data set obtained from a large DSP in the USA. Second, we evaluate the market level allocation algorithm and compare it to optimal offline methods and modified methods currently used to solve this class of matching problems. We use simulated data for the second evaluation since we do not have access to the full range of advertiser scores for each ad space in the DSP data set.

2.5.1 Single-Advertiser Simulation

To evaluate the agent level algorithm, we use real-world data from a DSP. We obtained a 14-hour window of ad space scores for a single campaign totaling 102,320 observations. In Figure 2-4, we plot the total ad spaces available for an advertiser using stacked score bands. From this chart, we can see that the total supply of ads is variable and that the score distribution over time is also changing. Figure 2-5 shows the same data in proportions instead to illustrate the non-stationarity of the series. Additionally, when we think of a thresholding rule for selecting the number of ad spaces, we can see the non-linearities in setting such a threshold. For example, we can see that adjusting the threshold in the 4-10 range has a much smaller effect on ads won when compared to adjusting the threshold in the 1-2 range. This is precisely the type of non-linearities that are ignored in many other mechanisms that assume known distributions and why a more complex controller like a PID controller is needed. These time series plots show just how variable each minute's supply is and that the underlying ad rate is non-stationary and why a time series smoother is necessary.

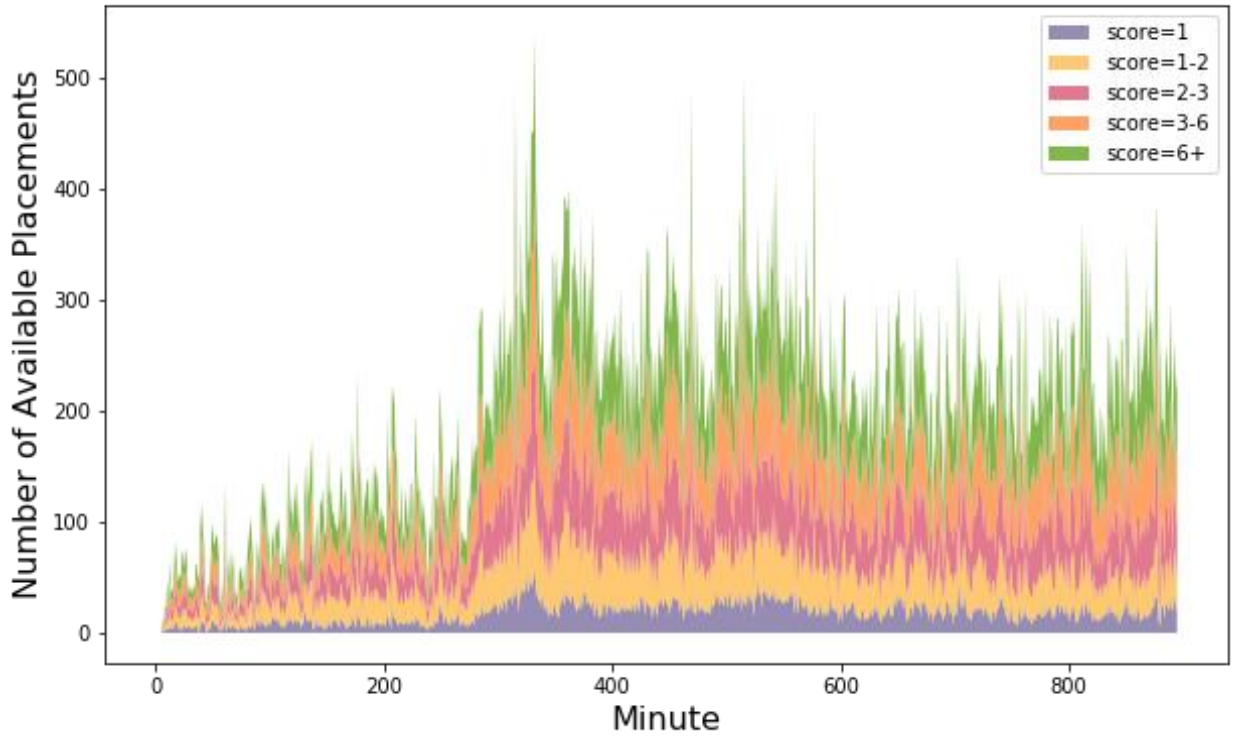


Figure 2-4 Score Distribution by Minute

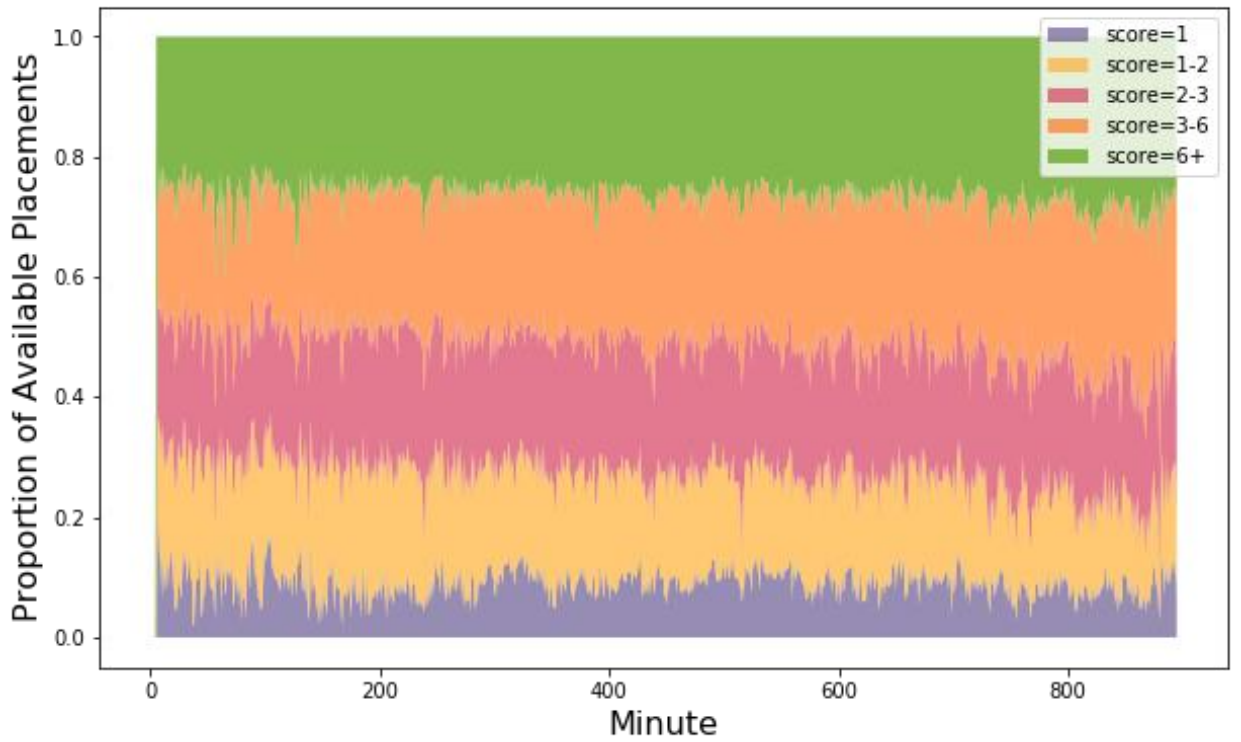


Figure 2-5 Normalized Score Distribution by Minute

We partition the data into four time periods: from minutes 0-299, 300-499, 500-699, and 700-840. The first time period is used for calibration purposes, our proposed method does not require this, but all competing methods do. In general, we make our evaluation choices in favor of the competing methods. Next, we set a constant minute level ad target (quota) for these four time periods: 30, 30, 5, and 45 ad spaces respectively. The changes in ad targets are equivalent to shocking the system to illustrate the ability of the competing methods to adjust to changing circumstance. The ideal outcome is that with perfect hindsight (i.e., an oracle condition) that we would select only the highest-scoring ad spaces needed in each time period to satisfy the quota. Our measures for success in the context of just a single advertiser are efficiency and delivery. Efficiency is calculated as the sum of the total ad scores chosen during that time period and then divided by the total sum that an oracle would have achieved. Thus, efficiency can be seen as a percentage of the theoretical maximum score possible. Delivery is calculated as the percentage of the quota that has been fulfilled based on the set ad target for that time period.

We consider two benchmark comparisons for the single advertiser version of our problem. First, we use the simple rule of buying all the ad space as it arrives until the quota is met (which we call First Ads). While trivial, this rule is used in practice given its simplicity and the importance of meeting the ad quota. Second, we retroactively look at the training period and pick a threshold with perfect hindsight that would yield the number of ad spaces required (which we call Baseline). This second approach is precisely what a primal-dual based solution (such as Feldman et al. 2018) would reduce to in a single advertiser setting. Our third set of results show the application of the proposed agent-level algorithm. We note that the agent-level algorithm is at a disadvantage in this simulation as it is the only one attempting to hit a minute-level impression target.

Table 2-5 Agent Level Evaluation Results

	Minute time interval		
	300-499	500-699	700-840
First Ads Efficiency	37%	18%	65%
First Ads Delivery	100%	100%	100%
Baseline Efficiency	74%	79%	90%
Baseline Delivery	178%	212%	201%
Our Approach Efficiency	97%	90%	95%
Our Approach Delivery	100%	100%	100%

Looking at the results in Table 2-3, it is unsurprising that the First Ads method always buys sufficient ads. Also, it is similarly unsurprising that efficiency is between 18% and 65% of the oracle and theoretical maximum. This method makes no effort to pick better ad spaces and focus entirely on meeting the quota.

The Baseline method does significantly better than the First Ad method. It happened to be calibrated in a time period where ad supply was lower than the later periods. So even though it was more selective with which ad spaces to select, it still selected too many. Hence, it over-delivered ads without an additional stopping rule once the quota was hit. The more selective thresholding improves efficiency, providing between 74% to 90% of the oracle’s performance.

Finally, our proposed, agent-level method adjusted dynamically in the first period (which we do not compare) and is well-calibrated at the start of the first test period. Given that it starts off optimally and adjusts quickly to the changing ad volume, it yields a 97% efficiency level relative to the oracle in the first test period. In the second and third test periods, it achieves 90% and 95%

efficiency respectively, marginally lower because of the time needed to adjust to the new requirements.

In Figures 2-6 and 2-7, we take a closer at the proposed method. Figure 2-6 shows the ad supply over time (which is not known by the algorithm), and the target ad quota. We begin with a threshold of 10, which quickly drops to 0 since the desired number of ads is initially above the supply. We also see the realized and smoothed number of ads obtained for each minute given the selected threshold. Once the supply exceeds demand, the threshold quickly rises, illustrating how the algorithm becomes rapidly more selective once better options are available. We also see that when the ad target suddenly drops and later rises that the threshold similarly rises and then drops rapidly to accommodate the new requirements.

In Figure 2-7, we plot the two tuning parameters of the PID controller to illustrate the learning process of the particle swarm optimization. The tuning parameters are initialized at a relatively slow-moving and stable point. Also, according to standard practice, the tuning parameters are bounded within a stable operating range of $0.005 < K_u < 0.1$ and $5 < P_u < 100$. At the start, the algorithm realizes that it is far from satisfactory and aggressively raises the K_u parameter. Once it can supply enough ads, the tuning parameters settle down into a less reactive state with low K_u and P_u values. Then later, when the ad targets suddenly change, it realizes that it needs to become drastically more reactive and manage much smaller quotas, so it raises both K_u and P_u values. This illustrates why a dynamic PID controller is required for the DSP setting. In Appendix B, we show additional results without using the PSO to demonstrate its necessity.

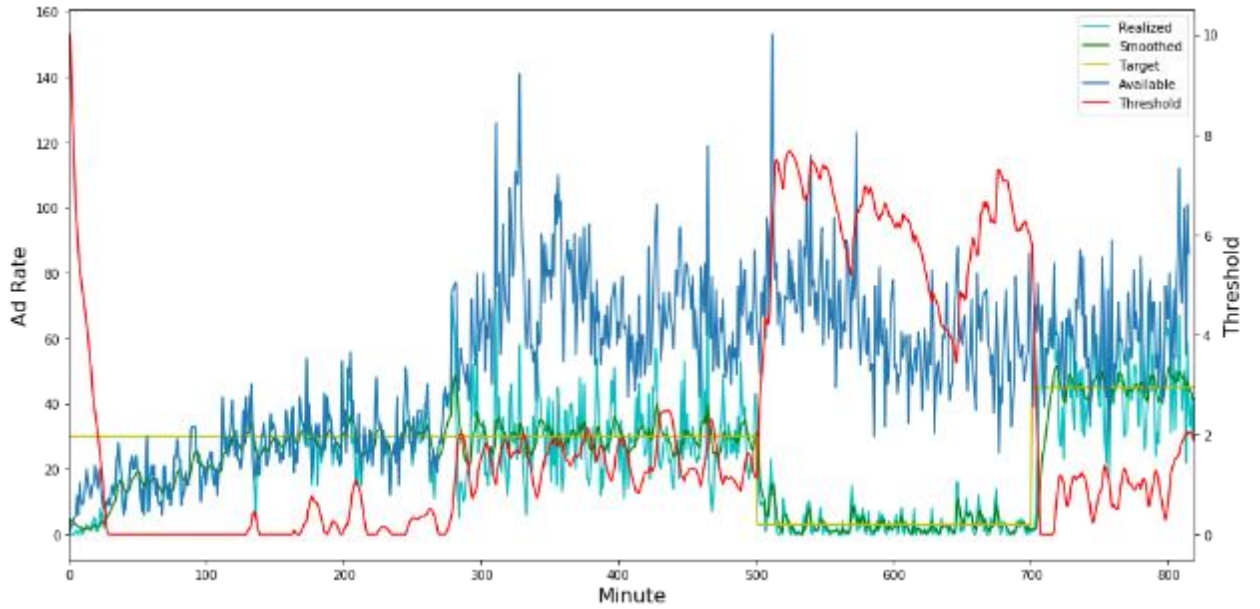


Figure 2-6 Simulation Results for Proposed Solution

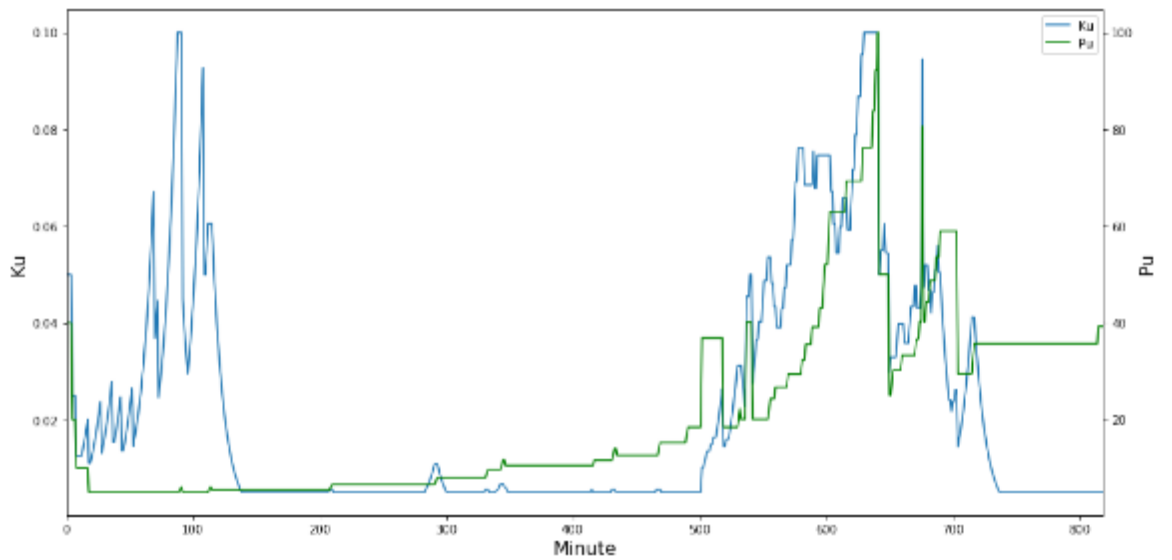


Figure 2-7 PID Parameters for Proposed Solution over Time

2.5.2 Multiple Advertisers (Market-level evaluation)

Next, we examine the market level algorithm, which is the mechanism that decides amongst all the advertisers that desire an ad spot, which one should get it. Unfortunately, we do

not have DSP data that reflects the scores of all advertisers for each ad space that appears on RTB exchanges. Therefore, we simulate data to approximate this environment. For each iteration, we generate preferences for three advertisers across five items. The item scores are drawn from a Uniform(0,1) distribution, then each item is perturbed by independent draws from a Normal(0,1/5) for each advertiser to create correlated advertiser preferences. The item frequencies are drawn from a Dirichlet(1/5) distribution to represent the high level of imbalance of online ad spots. Using this item-advertiser score matrix and item frequencies, we sample 100,000 items and allocate them according to the various methods. We then repeat this entire process 100 times to obtain median outcomes for competing methods.

We compare our algorithm to four other methods. First, we compare it to random assignment. Although very simple, it is the benchmark of truly fair assignment as each advertiser has an equal chance of receiving the desired item regardless of the score. The second method is a greedy assignment. Greedy assignment allocates the ad spot to the advertiser scoring it the highest. This is akin to assigning an ad spot to its highest value and will lead to the least equal outcomes. Both random assignment and greedy assignment are used in the literature and in practice due to their simplicity and for their fairness/efficiency arguments (e.g., Devanur and Hayes 2009, Manshadi et al. 2010, Gollapudi and Panigrahi 2014, Feldman et al. 2018). However, these methods represent the two extremes and can lead to undesirable outcomes that we hope to overcome with our method.

Given that, to our knowledge, no online mechanism exists to fairly allocate a stream of items in real-time, we compare our results to the ASM method discussed earlier (McAfee 1992). For this method, we wait until the entire set of 100,000 items has appeared. Then we cycle through each advertiser in turn and allocate to them their highest rated item of the remaining items. This

item is removed from the set and assigned to that advertiser. Then we move on to the next advertiser and repeat the process. This continues sequentially through the three advertisers allocating 1 item at a time until all the items have been allocated. While impossible for the DSP to implement in practice, this provides a good benchmark for both fair and efficient allocation of a set of items.

Next, we apply a primal-dual solution based on the Devanur and Hayes (2009) formulation, which we refer to as DH. Since this is a solution to the adwords problem, our implementation is slightly modified as we do not have an advertiser budget, we instead create a budget based on the observed advertiser ad space scores during the training (see Appendix B for the implementation details). We allocate the first 10% of ad spots (i.e., the first 10,000) randomly. Then we use this training set and a quadratic solver to solve the primal-dual problem. For future iterations, the advertiser weights are multiplied by the advertiser's score for the drawn item. The item is then allocated to the advertiser with the largest of the updated scores.

Finally, we replicate the DualBase solution discussed in Feldman et al. (2018) as the best in class solution to the publisher problem. We again allocate the first 10% of ad spots (i.e., the first 10,000) randomly. Then we use this training set and a quadratic solver to solve the primal-dual problem. For future iterations, the advertiser weights are subtracted from the advertisers' scores for the drawn item. The item is then allocated to the advertiser with the largest of the updated scores.

Our simulation setup provides the best-case scenario for both of the primal-dual solutions since their primary weakness is that they assume that the learned distribution will remain the same over time. This assumption is valid in this simulation, but it is not required by our method and does not hold in practice.

We consider two evaluation metrics. First, we again define efficiency as the sum of the scores for the advertiser that the item was allocated to. By this definition, the greedy algorithm will represent the upper bound of this metric. Second, to measure fairness, we use the multivariate Gini distance metric (Koshevoy and Mosler, 1997). This metric is minimized when each advertiser receives an equal number of each item and it is representative of the intrinsic fairness notion that each advertiser has equal rights to each item. We can use this metric in our simulation because we have designed the simulation such that each advertiser desires all items. The Gini equation is presented below where n is the number of advertisers, i and j are advertisers, and d is the number of items.

$$R_D(F_A) = \frac{1}{2dn^2} \sum_{j=1}^n \sum_{i=1}^n \left(\sum_{s=1}^d \frac{(a_{is} - a_{js})^2}{\bar{a}_s^2} \right)^{1/2} .$$

We note this is different from the Feldman et al. (2018) fairness measure where they assume that item scores are accurate and comparable, this allows them to define a fairness measure based on agents' item scores thereby creating a value-based fairness metric. However, in our setting, the scores are not directly comparable, making their measure ill-defined for our purposes. Gini distance works well for our purposes as it looks only at the distribution of items and not their scores, this makes it a good measure to capture the fairness differences between our competing methods.

In Figure 2-8, we plot the results of the simulation with Gini distance on the y-axis and efficiency on the x-axis. We plot the median results for all methods and normalize such that random takes on the value of 1 in both dimensions. We first note that, unsurprisingly, Random is the least efficient but the most equitable. Similarly, Greedy is the most efficient, being 28.7% more efficient than random, and also the least equitable with a Gini distance almost 248 times larger

than random. DH is approximately 8.5% more efficient than Random while achieving a Gini distance almost 218 times larger than random. The DualBase solution is very close to greedy and is 26.3% more efficient than random with a Gini distance 224 times larger than random. ASM is 7.6% more efficient than Random, with only 16.3 times higher Gini distance when compared to random.

Our proposed method is shown in red in Figure 2-8. The red dots illustrate the median outcomes of different choices of γ , namely: 1/10, 1/5, 1/2, 1, 2, 3, 5, 10, 20, and 50. For example, the outcome with $\gamma=1$ is close to the ASM method, being 6.4% more efficient than random while achieving a distance only 35.1 times higher than random. These results illustrate how we can achieve a sliding scale from random to greedy and form an efficiency-fairness frontier allowing a tradeoff between the two. We are able to achieve greater levels of efficiency than the DH and DualBase implementations with more equitable outcomes. We note that we are able to come remarkably close to the ASM method, but given that it is the only offline solution, it is also able to perform beyond our efficiency-fairness frontier. Additionally, since our method requires linear time, it could be a useful approximation to applications of ASM where the offline datasets are too large for ASM.

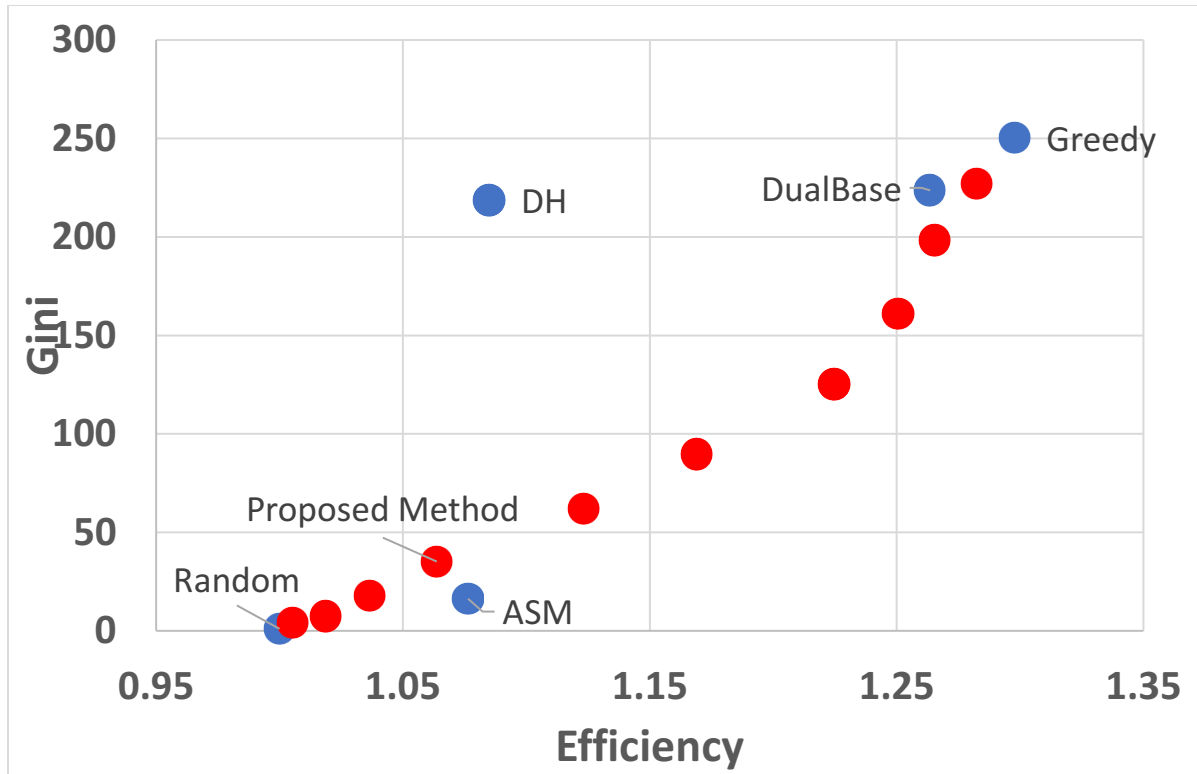


Figure 2-8 Simulation Results

2.6 Conclusion

Digital advertising continues to evolve rapidly, and the advent of real-time bidding is no exception. Now that publishers sell ad spaces from their websites in real-time, advertisers need to adapt to these new highspeed platforms. These changes have given rise to demand-side platforms who algorithmically purchase ad spots on behalf of advertisers and are tasked with fairly executing in the best interest of multiple advertisers simultaneously.

We present the DSP problem and show how it is different from the adwords and publisher problems in that ad supply is far more variable, ad quotas are more important, and fairness plays a larger role. This eliminates the possibility of relying heavily on earlier data for model training and heightens the need for a streaming algorithm that can adapt rapidly adapt to changing environments.

In this paper, we present such a streaming algorithm based on previous research in both economics and computer science. Through simulations using real-world ad exchange data and simulated data, we demonstrate that our proposed method is superior to current state-of-the-art methods for both the adwords and publisher problems as well as best-in-class offline methods.

Our proposed method exposes the explicit tradeoffs that DSPs need to make when assigning ad space between multiple advertisers that they represent. We show that there is also a firm-level tradeoff between fair ad allocations and efficient ad allocations. Our method provides a tuning parameter that allows DSPs to move between an entirely fair random allocation mechanism with low efficiency, to a greedy allocation mechanism with high efficiency but also highly unequal allocations. The parameter allows the DSP to pick a point on a frontier between these two extremes that can meet their business needs and is transparent to advertisers.

As a caveat, we offer no formal proof of optimality, in particular, we do not show that the tradeoff frontier we provide is Pareto optimal. It is challenging to provide such a rigorous proof for an algorithm that is dynamic and state-dependent. We leave this as an open question for future research.

We also acknowledge that fairness is a highly subjective issue. We attempt to address this by separating our efficiency measure from our fairness measure rather than finding a unified measure to assess the methods. We use the well-accepted Gini measure for fairness, which is both familiar to many but also treats all items equally in terms of value to fairness.

3 Cross-Merchant Spillovers in Coalition Loyalty Programs

3.1 Introduction

Loyalty programs have become a fundamental marketing tool for firms seeking to improve outcomes ranging from customer retention to data collection. In a typical loyalty program design, members collect points (e.g., one point per dollar spent) from purchases at a specific merchant, and spend these points on future purchases with the merchant (e.g., at a rate of a dollar per 100 points earned). Many highly-subscribed programs follow this basic model; well-known examples include Starbucks Rewards and CVS ExtraCare. Coalition loyalty programs are groups of firms within which consumers can earn and spend points at any merchant. These coalitions promise customers greater flexibility in the way they earn and redeem points with the goal of increasing overall participation within the network of merchants. Coalition loyalty programs attempt to provide greater value by increasing engagement, providing targeted offers using broader user data, and including non-competing partners at which to redeem points (Sports Loyalty International, 2018).

Well-known examples of coalition loyalty programs have been met with varying levels of success over the years. Germany's Payback coalition program, which began in 2000 and has expanded into five other countries, boasts 30 million active card users and six million active app users (Payback, 2018). Part of the American Express group, Payback enables its users to earn and use points at hundreds of merchants. Payback claims to have "generated" 33.8 billion euros in 2018, 409 million euros worth of points collected in 2018, and a 95 percent redemption rate of all points collected (Payback, 2018).

Another large coalition is the United Kingdom's Nectar, which was founded in 2002. As of 2016, Nectar had 19 million users and one million app downloads (Hobbs, 2016) and currently

has hundreds of participating merchants (Nectar, 2018). Other examples are Air Miles Canada and Travel Club in Spain (both reaching 70 percent household penetration), Fly Buys in Australia (60 percent household penetration; Sports Loyalty International, 2018), and South Korea's OK Cashbag (Yi et al. 2014).

A recent example of an unsuccessful coalition loyalty program is Plenti, a program in the U.S. that closed in July 2018 after three years in operation. Merchants in the Plenti coalition included firms with a national presence such as Macy's, a department store; Exxon and Mobil gas stations; Rite Aid, a pharmacy; and Chili's, a restaurant chain. Owned by American Express, Plenti had 36 million users at its peak (Pearson, 2018); however, very few were redeeming points for rewards, and those that did rarely made more than two redemptions (Shoulberg, 2018).

While loyalty program partnerships occur in many forms, we focus on coalition loyalty programs where a curated group of merchants share a common loyalty point pool. This is distinct from other partnership programs such as credit card points where the points are earned at any of a vast number of merchants accepting a payment method such as Visa or MasterCard. Credit card reward programs, moreover, are largely funded by credit card companies themselves as a means of driving adoption, whereas the coalition loyalty program we study is funded by individual merchants. The coalition program we study shares some similarities with airline alliances, in that points are easily transferrable between vendors and that vendors in the networks are generally non-competing. However, whereas airline alliances offer products in a single vertical, the loyalty program we study is designed to multi-category, presumably making the nature of spillovers quite different between the two cases.

The underlying thesis of coalition loyalty programs is that, relative to standalone programs, coalition programs more effectively enhance customer loyalty by allowing points or rewards to be

redeemable across a selection of merchants. Coalition programs also make possible lower customer acquisition costs between network merchants by facilitating cross-promotional activity. However, joining a coalition also exposes merchants to the risk of losing engagement to other network merchants; for instance, point redemption may be more attractive at some merchants than others.

We aim to contribute to the understanding of coalition loyalty programs by studying cross-vendor benefits from new merchant entry into a coalition loyalty program. Our analysis focuses on a coalition loyalty program in a Middle Eastern country that consumers use via a mobile application. We employ a novel identification strategy that relies on the entry of merchants into the coalition. These merchant entries comprise exogenous shocks that have, to our knowledge, not been studied in prior research. We use a dual strategy composed of regressions on matched samples and a Bayesian structural time series model to study consumers' responses to these events.

In our first approach, we use a regression framework to measure the impact of merchant entry on the pre-existing merchants in the network. We define the base group as the set of consumers who do not purchase at the new merchant; the relevant group as the set of consumers who purchase at the new merchant; and the dependent variables of interest as the outcomes for pre-existing merchants. Our key assumption is that merchant entry is irrelevant for users of the program who do not purchase at the entrant. We find that merchant entry has several positive effects on pre-existing merchants: more frequent transactions, larger basket sizes, and higher overall sales.

In our second approach, we use a Bayesian structural time series model to estimate the spillovers from merchant entry. The model forecasts the counterfactual paths for key outcomes had entry not occurred; the difference between the counterfactual and realized paths constitute the

treatment effect. The benefit of this approach is that it does not require the necessary assumptions required for our matching analysis. Nonetheless, we find convergent evidence of significant positive spillovers from merchant entry from our two approaches. In the following sections, we discuss the related literature, data and industry setting, our dual empirical strategies, and conclude.

3.2 Related literature

The bulk of the literature on loyalty programs is focused on standalone programs, i.e. those in which points are earned through purchases at one firm and rewards are redeemable at the same firm. In this section, we highlight key results from this literature, which range from findings on consumer decision-making to program design for firms. We then discuss the emerging literature on coalition loyalty programs.

Research in quantitative marketing has sought to understand how loyalty programs impact key consumer outcomes. For example, Lewis (2004) develops a dynamic structural model to study customer retention in a grocery loyalty program and finds that the program increases repeat-purchase rates. Similarly, Jiang, Nevskaya, and Thomadsen (2017) find a 15-17 percent reduction in customer attrition from a non-tiered loyalty program at a hair salon. Liu (2007) finds that customers who are initially heavy users in a loyalty program are more likely to redeem their rewards; however, moderate users tend to become more loyal to the focal store and increase both the number and size of their transactions over time.

In addition to assessing outcomes, related research has also evaluated how loyalty programs are structured. Kopalle, Sun, Neslin, Sun, and Swaminathan (2012) study the sales impact of a joint frequency and customer-tiered loyalty program and conclude that both components have significant effects on incremental sales. Zhang and Breugelmans (2012) find

several sources of gains in switching from a discount-based loyalty program to an item-based loyalty program. Wei and Xiao (2015) investigate the relative effectiveness of discounts and increased rewards points value, and find that unlike discounts, rewards promotions increase the number of purchases of other products in similar product categories. Using a dynamic structural model and data from an Italian gasoline merchant, Rossi (2017) finds that most consumers prefer price discounts instead of receiving rewards points and are insensitive to changing reward structures.

The interaction of promotions and loyalty programs, in particular, has been an emphasis on related research. Van Heerde and Bijmolt (2005) study the responsiveness of different customer groups using a Bayesian hierarchical framework and conclude that non-targeted price promotions are more profitable than techniques that target loyalty program members. Relatedly, Zhang and Wedel (2009) develop a model for optimizing campaigns that leads to significant profit bumps, as well as evaluating the effectiveness of different types of promotions for a combined online and offline retailer.

Some papers have focused on consumer behaviors that are idiosyncratic to loyalty programs. Stourm, Bradlow, and Fader (2015) develop a structural, multi-account model of a linear loyalty program and find three distinct causes for stockpiling of points: an economic incentive; a cognitive incentive; and a psychological incentive, with the latter two being dominant forces in their empirical setting. Orhun and Guo (2018) find that airline loyalty members are more likely to sacrifice near-term utility to gain status in tiered programs when they are “behind schedule” in their progress.

Researchers have also studied how loyalty programs interact with consumer psychology to influence behavior. Dreze and Nunes (2004) find that consumers prefer to use a combination of

cash and rewards points when they do not value cash and rewards points equally and when the perceived cost function for one of the currencies is partly convex. Meanwhile, using multiple surveys, Kwong, Soman, and Ho (2011) find that consumers are more likely to spend their points when the benefits to doing so are easily calculated, e.g. an easy to compute fraction.

There is a small but growing literature on coalition loyalty programs. Dorotic, Fok, Verhoef, and Bijmolt (2011) do not find evidence of cross-vendor spillovers from promotions in coalition loyalty programs. Danaher, Sajtos, and Danaher (2017) model the transition of consumers between behavior states in the context of a coalition loyalty program. Stourm, Bradlow, and Fader (2018) find that store affinity, e.g. similarity of product categories sold or geographic proximity, yields positive cross-reward effects, except in situations where partners are similar along both these dimensions (i.e. competitors are close to each other). To our knowledge, prior research has not examined the impact of network size or composition on the effectiveness of coalition loyalty programs. We aim to contribute to the literature by examining the impact on firms in a coalition loyalty program of the entry of new merchants into the network.

3.3 Background, data, and preliminary evidence

The data used in this study was provided by a coalition loyalty program in a Middle Eastern country that consumers use via a mobile application. The program operates in an economy that is experiencing improving internet connectivity, significant smartphone adoption, and significant mobile phone usage for shopping, social media, mobile banking, and activity engagement. Furthermore, the country's retail sector has experienced steadily increasing growth and is projected to continue growing for several years.

The coalition loyalty program operates by offering promotional campaigns at partnering merchants via its mobile applications. The program includes merchants in several industries including traditional retail, travel, mobile carrier, electronics, and grocery among others, and offers three types of campaigns: campaigns where customers can earn up to a specific number of points, campaigns where customers earn points as a percent of their total basket size, and campaigns where customers can redeem points at a higher value than the base exchange rate to receive greater discounts (hereafter multipliers).¹⁶ The loyalty program earns 77 percent of its revenue from commission when program members utilize promotional campaigns. The remaining revenue comes from point redemption, expiration, and additional advertising.

The data set provided to us contains detailed information on consumer activity and tracks unique users throughout the data set. For example, for each transaction we have the original basket size of transactions prior to any rewards points and multipliers being used. We also know the rewards points being earned by consumers and can identify the associated campaigns being utilized. The centerpiece of our research is an analysis of the entry of a large grocery chain into the coalition. We only consider the transaction data for the 3 months preceding the merchant (i.e. the beginning of our data set) and 3 months after entry. We constrain data for this analysis to maintain a stable set of shoppers, as well as to prevent contamination from further merchant entry. More detail for the full and subsample data sets and distributions of relevant variables are provided in Tables 3-1 and 3-2, respectively, below.

¹⁶ Program members also earn points at a rate of 0.1 percent of total basket size of all non-campaign transactions.

Table 3-1 Data Set Summary Statistics

	<u>Entire time period</u>		<u>Analysis time period</u>	
	June 3, 2015 - August 25, 2016		August 8, 2015 - February 4, 2016	
Duration				
Number of unique customers	1,538,638		681,777	
Number of unique transactions	8,274,819		2,336,141	
Number of unique merchants	73		47	
Number of unique stores	3,532		2,477	

Table 3-2 Distribution of Relevant Variables at the Transaction Level

	<u>Entire time period</u>		<u>Analysis time period</u>	
	Mean	Std. dev.	Mean	Std. dev.
Original basket sizes	139.57	132.824	168.02	141.77
Points Earned	32.46	27.25	37.98	29.23
Points Used/Redeemed	35.2	30.81	40.85	31.95

Sample means are statistically different from each other at $p < 0.001$ level

In our empirical analysis, we exploit the entry of merchants into the program as a source of exogenous shocks to customer behavior at pre-existing merchants. We plot the cumulative entry of merchants and their stores over time in Figures 3-1 and 3-2 below. As is evident from Figure 3-2, the merchant entry introduced a significant number of stores into the loyalty program, accounting for approximately one-half of all stores included in the analysis. In subsequent analyses, we also study the impact of smaller merchant entries.

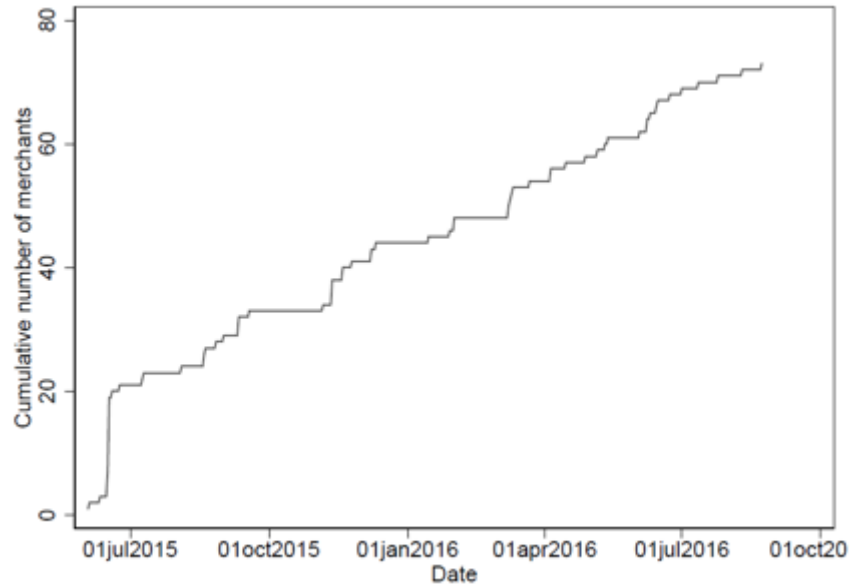


Figure 3-1 Cumulative number of merchants in the coalition

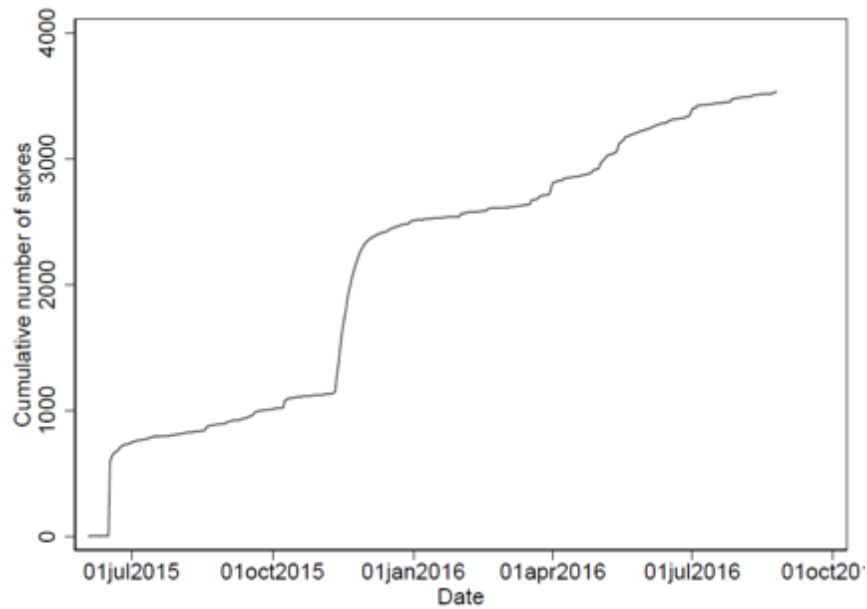


Figure 3-2 Cumulative number of stores in the coalition

In fact, we can see model-free evidence of positive externalities from more merchants and stores entering a coalition loyalty program. The analysis is performed at the store-day level and

does not include the large merchant that enters the program. Table 3-3 presents results from regressions on our main dependent variables: the total number of transactions across all merchants, average basket size for purchases, aggregate sales, rewards points earned, and rewards points redeemed. We control for network size (measured as the cumulative number of merchants in the program), distance to the nearest store, and a time trend. We include distance to nearest store as an additional measure of network growth; because this is mostly offline retail an in-network store that opens nearby potentially provides an additional shock dependent variables. These regressions provide exploratory evidence that increasing the size of the network leads to increases in each of our dependent variables for every store in the network on a daily basis. However, the distance to nearest store and time trends are sometimes negative and sometimes positive indicating that net spillovers may also be negative or positive and thus a more careful analysis is needed.

Table 3-3 Preliminary evidence with distance to the nearest store

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (000s)	(4) Points Earned	(5) Points Used
Merchant Network size	0.2675*** (0.0372)	17.96*** (1.01)	0.2158*** (0.0160)	20.32*** (2.614)	10.31*** (1.620)
Distance to Nearest Store	0.0226*** (0.0045)	-1.493*** (0.393)	0.0046** (0.0023)	3.677*** (0.285)	4.377*** (0.182)
Time Trend	- 0.0094*** (0.0036)	1.594*** (0.089)	0.0094*** (0.0014)	0.099 (0.269)	1.134*** (0.174)
Observations	172,058	172,058	172,058	172,058	172,058
R-squared	0.63	0.23	0.75	0.19	0.26

Day-of-week FE	YES	YES	YES	YES	YES
Store FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

In the following sections, we describe our approach for detecting the presence of spillovers from the entry of merchants into a coalition loyalty program. We adopt a two-pronged measurement strategy. The first consists of regressions on matched samples, where groups are selected based on eventual activity at the entering merchant. The second consists of counterfactual simulations using Bayesian structural time series modeling. These two distinct approaches produce convergent evidence of significant positive spillovers from merchant entry in the coalition loyalty program.

3.4 Matching Analysis

To estimate the impact of the entry of the large grocery chain in our sample, we partition our data set into two groups: a group of consumers that never shop at the new large merchant and a group of consumers that we observe to make at least one purchase at the new merchant. We call these two groups our *base* and *relevant* groups, respectively.¹⁷ To mitigate potential bias in our estimation, we restrict our data set to consumers that made at least one purchase (at any store) in

¹⁷ Our data set provides very granular geographic detail on stores (i.e. their longitudes and latitudes). Using this information, we could, in theory, categorize stores based on their distance to the nearest store of the new merchant that entered the coalition loyalty program and use a difference-in-differences approach. However, upon the new merchant's entry, over 1,200 stores enter the coalition program, and more than 95 percent of pre-existing stores are within two miles of the new merchant.

the 90 days prior to the new merchant entry; doing so allows us to study how active consumers' behavior changes without the results being confounded by new customers that enter the program because of the new merchant.

We balance our base and relevant groups on variables that reflect their activity in the program prior to the merchant's entry. Specifically, we balance the sample of customers based on their aggregate sales, average basket size, the number of unique stores they shop at, and the number of days on which they transact in the 90 days prior to the grocer's entry.

Table 3-4 displays a basic two-sample t-test to highlight the existence of imbalance between the two customer groups. Each variable, with the exception of aggregate sales, is unbalanced between our treatment and control group, indicating that matching consumers between the two groups will allow us to draw more robust conclusions from our ensuing regression analysis.

Table 3-4 Summary statistics before matching

Variable		Control	Treated	Difference	
				C - T	P-value
Aggregate sales	Mean	1,137.46	1,170.89	-33.43	0.231
	Std. Error	8.89	24.58	27.92	
Average basket size	Mean	501.63	430.71	70.92	0.000
	Std. Error	2.58	6.24	8.03	
Number of stores visited	Mean	1.57	1.83	-0.26	0.000
	Std. Error	0.0025	0.0095	0.0081	
Number of days making a transaction	Mean	1.99	2.44	-0.45	0.000
	Std. Error	0.0043	0.0164	0.0139	
Number of observations		144,908	16,058		

Table 3-5 displays the results of the logistic regression performed to estimate the propensity scores. We follow the standard procedure described in Rubin (1974), where the binary independent variable in our logistic regression has a value of 1 if a consumer is observed to purchase at the grocery store and 0 otherwise, the number of observations represent the number of consumers. By

matching our samples, we aim to form two similar groups of consumers for comparison, and to mitigate any bias arising from the selection of one group into purchasing at the grocery store. Table 3-6 presents summary statistics of the matching variables of the propensity score matched sample against the unmatched sample. Despite the relatively low pseudo-R² scores, our propensity score matching still manages to achieve significant reductions in imbalance across all variables indicating that it does explain a large portion of purchasing behavior.

Table 3-5 Propensity score logit regression, matching

Variable	Treated Customer
Aggregate sales	-0.00003*** (0.00001)
Average basket size	-0.000101*** (0.00002)
Number of stores visited	0.14032*** (0.00995)
Number of days making a transaction	0.09660*** (0.00704)
Constant	-2.58178*** (0.00678)
Observations	160,966
Log Likelihood	-51,597.824
Pseudo R-squared	0.0123

Table 3-6 Summary statistics after matching

Variable		Control	Treatment	Difference	
				C - T	P-value
Aggregate sales	Mean	1170.9	1208.8	-37.9	0.277
Average basket size	Mean	430.71	436.5	-5.79	0.507
Number of stores visited	Mean	1.8294	1.8434	-0.014	0.297
Number of days making a transaction	Mean	2.4364	2.4342	0.0022	0.926
Number of observations		16,058	16,058		

We focus on outcome variables specific to incumbent merchants, i.e. excluding activity at the new merchant. Our outcome variables of interest are: the overall number of transactions; average basket sizes (net of all discounts, i.e. points used plus applied multipliers); aggregate sales (also net of all discounts); points earned; and points used (see Appendix C for variable definitions). We estimate the following program-level model:

$$y_{i,t} = \alpha + \beta_1 Post\ Entry + \beta_2 Customer\ Group_i + \beta_3 Post\ Entry * Customer\ Group_i + \theta Fixed\ Effects + \varepsilon_{i,t} \quad (1)$$

where $Customer\ Group_i$ indicates whether the customer falls into the base group or relevant group (i.e. grocery store customers), $Post\ Entry$ is an indicator for pre- and post-merchant entry, t is the day in the 181-day period analyzed, and $Fixed\ Effects$ is a matrix of additional controls. Our main results, aggregated to the daily level for the program, are presented in Table 3-7.

The effect under investigation is post-entry impact on grocery customers for existing coalition merchants (i.e. the Table 3-7 parameter of Post-entry * Grocery customers). For this group we find that the total number of transactions increases by about 153 per day. Next, we note that the average basket size of each transaction is 43 units of local currency lower but that the aggregate sales is approximately 37,962 units of local currency higher. We also see that this group earns and spends approximately 3,454 and 3,090 loyalty points respectively.

Next, we perform a similar analysis to determine the effect of the new merchant entry at the individual store level for all pre-existing and non-grocery store merchants. We add the store level analysis to investigate if the aggregate analysis may be masking compositional effects, for

instance perhaps only the largest merchants experience positive spillovers. We adapt the model above as follows:

$$y_{s,i,t} = \alpha + \beta_1 \textit{Post Entry} + \beta_2 \textit{Customer Group}_i + \beta_3 \textit{Post Entry} * \textit{Customer Group}_i + \theta \textit{Fixed Effects} + \varepsilon_{s,i,t} \quad (2)$$

where i and t are defined as before, and s represents each store. We arrive at the results listed in Table 3-8.

Table 3-7 Aggregate level

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (000s)	(4) Points Earned	(5) Points Used
Post-Grocery entry	-6.37 (17.07)	22.71** (11.1593)	12.07 (8.77)	5,114.04*** (1,267.51)	6,767.25*** (882.85)
Grocery customers	1.87 (20.09)	-18.42* (10.77)	-6.75 (9.06)	1,148.70 (1,402.24)	145.22 (634.62)
Post-entry * Grocery customers	153.08*** (27.70)	-43.05*** (15.06)	37.96*** (14.31)	3,454.37* (1,942.80)	3,090.07** (1,362.73)
Constant	522.91*** (13.32)	418.97*** (7.43)	215.63*** (5.75)	7,800.40*** (917.69)	6,796.62*** (410.91)
Observations	362	362	362	362	362
R-squared	0.52	0.13	0.38	0.27	0.41
Day-of-week FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

Table 3-8 Store level

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (000s)	(4) Points Earned	(5) Points Used
Post-Grocery entry	-0.04*** (0.01)	-15.60*** (2.56)	0.014 (0.009)	10.79*** (1.04)	13.82*** (0.91)
Grocery customers	-0.003 (0.01)	-14.75*** (2.78)	-0.030*** (0.009)	1.54** (0.71)	-0.25 (0.52)
Post-entry * Grocery customers	0.25*** (0.02)	14.84*** (3.44)	0.063*** (0.013)	5.67*** (1.55)	5.67*** (1.27)
Constant	1.34*** (0.02)	170.84*** (2.80)	0.519*** (0.010)	26.17*** (1.39)	21.40*** (1.05)
Observations	168,037	168,037	168,037	168,037	168,037
R-squared	0.69	0.25	0.57	0.09	0.12
Day-of-week FE	YES	YES	YES	YES	YES
Store FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

We again focus on the post-entry impact on grocery customers (i.e. the parameter of Post-entry * Grocery customers). Here we see that that average store gets 0.25 more transactions per day. The average basket size now increases by 14.8 units of the local currency and the total sales increases by 63 units of the local currency. Again, both loyalty points earned and spent increases by 26 and 21 points, respectively.

From Table's 3-7 and 3-8 we draw our main conclusions, that stores see statistically significant increases in the number of transactions and aggregate sales from relevant customers at pre-existing coalition stores. Similarly, the relevant customers become more active within the loyalty program as can be seen by the statistically significant increases loyalty points earned and spent. The effect on baskets size implies that, while basket sizes increase at the store level, the estimated increase in the number of transactions accrues disproportionately to stores at which consumers have lower basket sizes.

Next, we seek to show how spillovers vary according to the intensity at which relevant customers transact with the entering merchant. We partition the set of relevant customers into terciles according to their total spending at the entering merchant. We add interactions to our baseline models indicating whether customers are low, medium, or high value customers according to this measure. Tables 3-9 and 3-10 contain the corresponding regression results.

We again focus on the post-grocery entry group interacted with the three tiers of customers. Across all tiers of customers we again see that total transactions increases and that the average basket size decreases. The pattern seems strongest for the lowest value customers who may coincidentally be the most responsive to promotional activity. However, we find no effect on aggregate sales, mixed results on points earned, and a decrease in points spent¹⁸.

¹⁸ This could be due to the lower power from the reduced number of customers in each tercile.

Overall, we find that after the grocery merchant joins the coalition loyalty program, the pre-existing customers who shop at the grocery store increase their number of transactions and money spent at other coalition merchant stores. In appendix C, for the sake of robustness, we repeat the analysis of this section using the unmatched data and report qualitatively similar results.

Table 3-9 Heterogeneous effects: Aggregate level

VARIABLES	(1) Total # of transactions	(2) Average basket size	(3) Aggregate sales (000s)	(4) Points Earned	(5) Points Used
Post-Grocery entry	-6.58 (19.31)	22.72** (11.25)	12.04 (9.45)	5,110.50*** (1,288.80)	6,761.75*** (906.50)
Low-value Grocery customers	-358.36*** (15.39)	-46.78*** (12.93)	-154.95*** (6.48)	-4,686.80*** (1,000.25)	-4,573.16*** (429.19)
Mid-value Grocery customers	-174.40*** (7.81)	-9.07 (6.48)	-73.15*** (3.32)	-2,502.28*** (493.53)	-2,235.12*** (213.08)
High-value Grocery customers	-112.42*** (5.26)	1.42 (3.87)	-45.60*** (2.24)	-1,588.70*** (330.89)	-1,473.07*** (147.10)
Post-Grocery entry *	53.11**	-48.28***	1.63	-2,193.85	-3,288.09***
Low-value Grocery customers	(20.74)	(17.42)	(10.24)	(1,429.45)	(1,022.07)
Post-Grocery entry *	25.93**	-23.40***	1.11	-1,404.77**	-2,012.53***
Mid-value Grocery customers	(10.50)	(8.85)	(5.18)	(685.65)	(482.77)
Post-Grocery entry *	20.57***	-10.46*	3.35	-585.58	-1,035.69***
High-value Grocery customers	(7.12)	(5.50)	(3.55)	(466.06)	(330.27)
Constant	523.04*** (14.42)	418.97*** (7.46)	215.65*** (5.92)	7,801.70*** (919.72)	6,800.20*** (381.06)
Observations	724	724	724	724	724
R-squared	0.784	0.130	0.738	0.333	0.502
Day-of-week FE	YES	YES	YES	YES	YES

Table 3-10 Heterogeneous effects: Store level

VARIABLES	(1) Total # of transactions	(2) Average basket size	(3) Aggregate sales (000s)	(4) Points Earned	(5) Points Used
Post-Grocery entry	-0.127*** (0.015)	-20.840*** (2.598)	-0.020** (0.010)	8.957*** (1.011)	11.893*** (0.888)
Low-value Grocery customers	-0.798*** (0.013)	-90.528*** (2.618)	-0.352*** (0.007)	-10.919*** (0.500)	-10.594*** (0.372)
Mid-value Grocery customers	-0.384*** (0.006)	-43.483*** (1.329)	-0.161*** (0.004)	-5.415*** (0.247)	-5.013*** (0.196)
High-value Grocery customers	-0.249*** (0.004)	-27.813*** (0.880)	-0.101*** (0.002)	-3.521*** (0.164)	-3.324*** (0.136)
Post-Grocery entry *	0.223***	24.871***	0.057***	-2.135*	-3.940***
Low-value Grocery customers	(0.017)	(3.217)	(0.011)	(1.180)	(1.025)
Post-Grocery entry *	0.110***	12.896***	0.028***	-1.877***	-2.827***
Mid-value Grocery customers	(0.008)	(1.631)	(0.005)	(0.551)	(0.475)
Post-Grocery entry *	0.076***	8.885***	0.021***	-0.713*	-1.350***
High-value Grocery customers	(0.006)	(1.079)	(0.003)	(0.383)	(0.337)
Constant	1.260*** (0.013)	167.571*** (2.451)	0.501*** (0.008)	21.805*** (0.752)	18.483*** (0.589)
Observations	328,507	328,507	328,507	328,507	328,507
R-squared	0.590	0.198	0.384	0.063	0.076
Day-of-week FE	YES	YES	YES	YES	YES
Store FE	YES	YES	YES	YES	YES

3.5 Bayesian Structural Time Series

Our preceding analysis relies on finding a representative base group that is unaffected by the market intervention. Although our matching approach mitigates selection bias, the exogenous merchant entry was a national event, and we have no general subset of the population unexposed to this event. Our approach also uses a linear model to account for other sources of variation in each group over time. This ignores the temporal nature of the data set and may incorrectly assume that error components are i.i.d. Standard linear models also pose modeling difficulties for variable selection and model validation. Our preceding estimators assume that there is a level effect caused by the intervention and that the difference between the two groups represents the causal effect of the treatment (Bertrand et al. 2004, Hansen 2007). The effect of merchant entry may not constitute pure level effects and instead evolve over time. Regression methods are insufficient for providing an estimate of the duration effect of the marketing intervention.

To overcome these limitations, we estimate the effect of the new merchant entry on the existing network of stores using a Bayesian Structural Time Series (BSTS) model proposed by Scott and Varian (2014). This method generalizes our preceding approach to a time series setting by using a state-space model with a regression component. Using this approach, we estimate a time series model based on the pre-merchant entry data and predict a counterfactual time series for the post-merchant entry series. There are several advantages to this approach. The first is that the state-space model accounts for the serial correlation in the data. Second, it provides a fully Bayesian framework for variable selection. Third, it can estimate the evolution of the intervention effect over time and thereby provides a fuller estimate of the intervention effect's wear-in and wear-out.

The primary downside to the BSTS approach is that it assumes that we can adequately generate a counterfactual estimate of the dependent variables. This relies first on having enough independent variables and observations to adequately model pre-intervention dependent variables. Second, we need to assume that the post-intervention environment and the parameter link between the independent and dependent variables are unchanged.

In our setting, we estimate the counterfactual on five months of data with a wide variety of independent variables (see Appendix C for variables). After a first estimation, we estimate a final model by removing all variables without an inclusion probability of at least 80% to minimize multicollinearity concerns. We then extrapolate our counterfactual onto the following three months of data. We assume that the aggregate market relationship between promotional, seasonal, and other independent variables is unlikely to change substantially within a three-month period. Additionally, since the average and median time between consumer purchases are 32 and 10 days, respectively, we are still likely to capture any meaningful changes in consumer behavior. As in our preceding approach, we consider the subset of users that existed in the data set prior to the large grocer's entry who also eventually shopped there post-entry.

BSTS Specification

We specify the model in line with Brodersen et al. (2015) and Scott and Varian (2014) where we have a hidden state α linked to the measured time series of our outcome variables $\{y_t\}$ as a generalization of the classic constant trend regression model:

$$y_t = \mu_t + z_t + v_t \quad v_t \sim N(0, V) \quad (3)$$

Where μ_t is a local linear trend component:

$$\mu_t = \mu_{t-1} + b_{t-1} + h_t + w_{1,t} \quad w_{1,t} \sim N(0, W_1) \quad (4)$$

with slope equation:

$$b_t = b_{t-1} + w_{2,t} \quad w_{2,t} \sim N(0, W_2) \quad (5)$$

And a weekly seasonal component:

$$h_{t+1} = - \sum_{s=0}^{50} h_{t-s} + w_{3,t} \quad w_{3,t} \sim N(0, W_3) \quad (6)$$

z_t is a static regression component:

$$z_t = \beta^T x_t \quad (7)$$

All the error terms $\{v_t, w_{1t}, w_{2t}, w_{3t}\}$ are normally distributed with respective covariance matrices $\{V, W_1, W_2, W_3\}$.

We use the Bayesian form to specify the model as this is most natural for state-space models and it allows for the spike-and-slab prior for variable selection in the regression component (George and McCulloch, 1994; Madigan and Raftery, 1994). The spike-and-slab prior is a combination of a point mass (the spike) on each parameter being 0 and a diffuse component (the slab) giving weight to a wide range of values for each parameter. Here, γ is the binary (spike) parameter of whether or not a given β is 0 or 1 and is modeled as a Bernoulli random variable with parameter π :

$$\gamma \sim \prod_i \pi_i^{\gamma_i} (1 - \pi_i)^{1-\gamma_i} \quad w_{3,t} \sim N(0, W_3) \quad (8)$$

$$\beta|\gamma, \sigma^{-2} \sim N(b_\gamma, \sigma^2(\Omega_\gamma)^{-1}) \quad \frac{1}{\sigma^2} \sim \Gamma\left(\frac{df}{2}, \frac{ss}{2}\right) \quad (9)$$

$$\Omega_\gamma^{-1} = \frac{\kappa(X^T X)}{n} \quad (10)$$

The slab component is normally distributed conditional on γ with the standard conjugate Gamma variance component. The regression parameter b_γ is conditioned on not being 0. The term corresponding to Zellner's g-prior (Zellner, 1986), κ , is the effective number of observations worth of information. The inverse gamma variance is the standard form (Gelman et al., 2002) where df is the degrees of freedom and ss is the total sum of squares. For further details on the BSTS specification and computation see Brodersen et al. (2015) and Scott and Varian (2014).

Inference

We can estimate this model using standard Markov Chain Monte Carlo (MCMC) simulations to obtain the posterior distributions for the state vector \mathbf{a} and the model parameters collectively termed $\boldsymbol{\theta}$. Then we simulate draws from \mathbf{a} and $\boldsymbol{\theta}$ given the observed data $y_{1..n}$. Using the simulated parameters, we simulate the counterfactual $\tilde{y}_{n+1..m}$ observations from the posterior predictive distribution $p(\tilde{y}_{n+1..m}|y_{1..n}, x_{1..m})$. We can use the observed $y_{n+1..m}$ and simulated $\tilde{y}_{n+1..m}$ data to estimate the point-wise causal effect of the intervention for draw τ as:

$$\phi_t^{(\tau)} := y_t - \tilde{y}_t \quad (11)$$

The average daily causal effect is estimated by taking the mean over the $n+1$ to m point-wise estimates:

$$\frac{1}{m-n} \sum_{t'=n+1}^t \phi_{t'}^{(\tau)} \quad \forall t = n+1, \dots, m \quad (12)$$

Finally, by repeating this sampling process, we can estimate the mean point-wise causal effect and our credible intervals for the causal effect.

Results

Since we are interested in constructing a predictive model for our measurement variables, we are less concerned about the variables included in the z_t component and more concerned about their predictive ability. As such, we include all other non-measurement variables in the regression component and rely on Bayesian variable selection to keep only the relevant variables for prediction. For all the variables in this analysis, we again focus on purchase behavior at merchants other than the large grocer entering the program.

We follow the Bayesian variable selection process with model averaging recommended by George and McCulloch (1997). Here we fit many potential models using the spike-and-slab prior based on random draws from the parameter posterior distributions, then estimate the posterior probability that each is the best model, and finally average over those models based on these posterior probabilities. We do this by setting the parameters $\{V, W_1, W_2, W_3\}$ to $1/0.05$ as

recommended by Brodersen et al. (2015). We then estimate the full model across 1,000 draws which allows us to construct 1,000 counterfactuals and estimate the average causal effect and the variance of this effect along with the Bayesian credible intervals.

For the remainder of this section, all of the BSTS graphs can be interpreted in the same way. The data is plotted at the daily level and we use a vertical dotted line on November 6th to illustrate when the grocery store entered the coalition loyalty program. The results are presented in sets of three plots: original, pointwise, and cumulative. The original plot uses a solid time series line to represent the observed outcome data and the dotted time series line represents the predicted outcome. The BSTS procedure uses only the outcome data observed before November 6th to generate the predicted time series. This plot allows us to see how well the model fits the time series before the intervention as well as deviation from the counterfactual prediction after the intervention. The pointwise plot shows the daily difference between the observed and counterfactual data and is an illustration of how the intervention effect may change over time. The cumulative plot shows the running-sum of the pointwise plot illustrating the total effect of the intervention for the post-entry period. Finally, the shaded area represents the 95% Bayesian credible interval obtained from the 1,000 MCMC samples.

Our first measure is the average basket size of purchases at pre-existing stores (see Figure 3-3). We note that this series increases towards the end of the year shopping season, and decreases later in the observation period. In spite of the large basket spikes around the holiday and New Year's shopping period, the observed basket size shrinks well below the predictions of the counterfactual. We estimate that the average basket size decreases by 54¹⁹ with a 95% Bayesian

¹⁹ Currency units are in the local currency, omitted for privacy purposes.

credible interval of [-175, -2]. The posterior probability of a positive causal effect is 96%. This is equivalent to a 12% decrease in the observed basket size with a 95% Bayesian credible interval of [-40%, -0.5%].

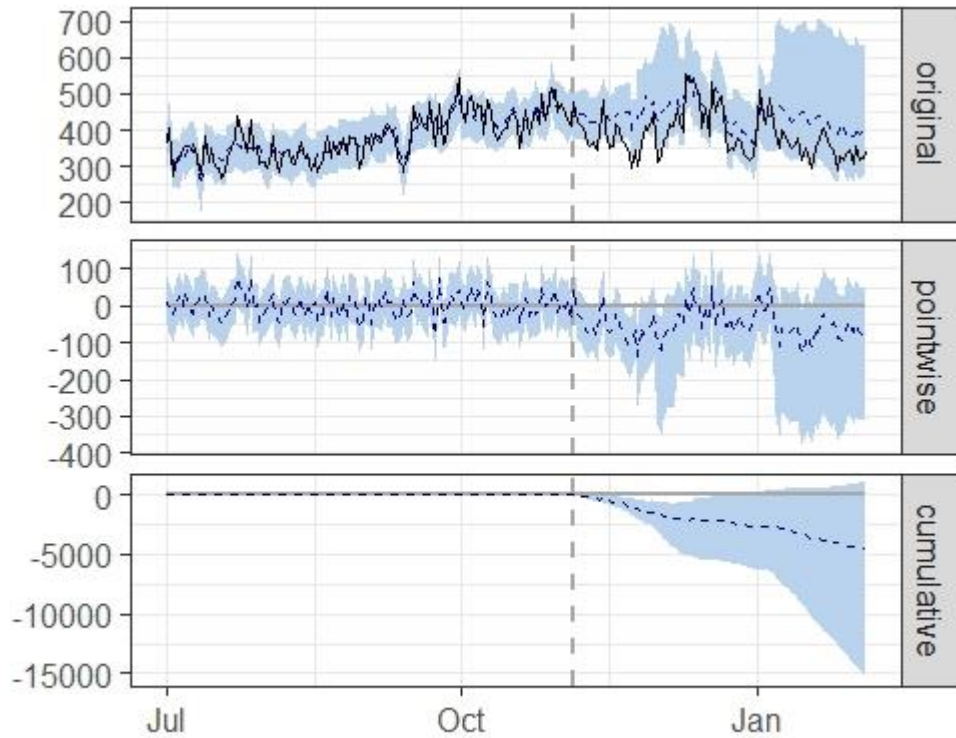


Figure 3-3 Comparing observed data to counterfactual data for average basket size

Our second measure is the total value of purchases at pre-existing stores per day (see Figure 3-4). We note that this series increases steadily throughout the observation period, peaks in the holiday season and dips in January. We estimate that total daily sales for the relevant cohort at pre-existing stores increase by around 150,000 local currency with a 95% Bayesian credible interval of [40,000, 250,000]. The posterior probability of a positive causal effect is 99%. This is equivalent to a 9.5% increase in daily sales with a 95% Bayesian credible interval of [7%, 44%].

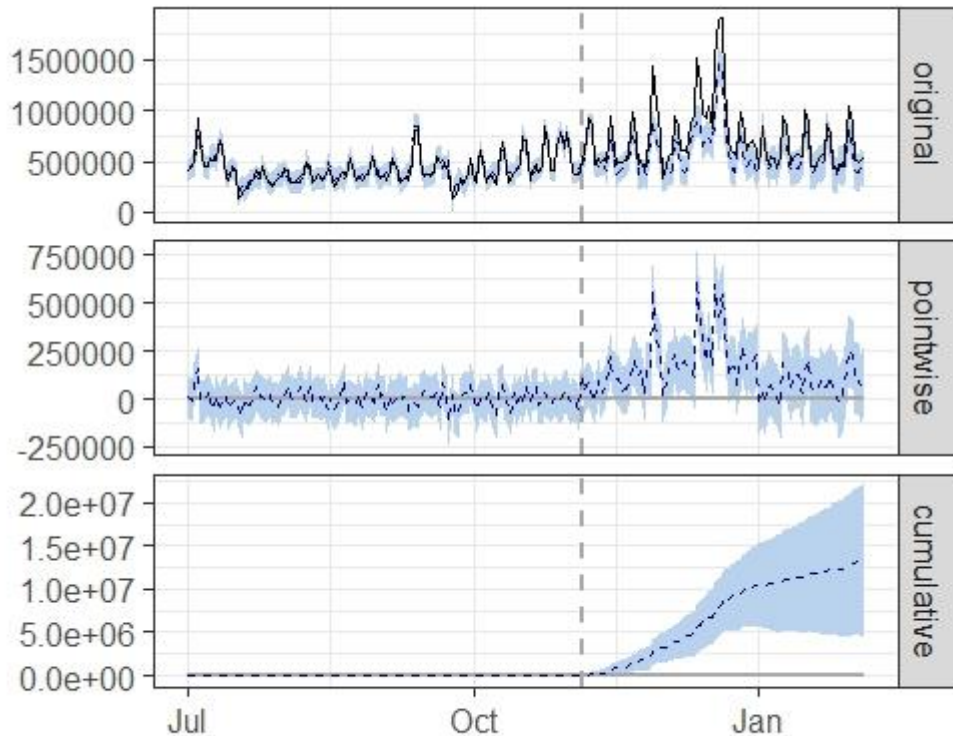


Figure 3-4 Comparing observed data to counterfactual data for aggregate sales

Our third measure is the total number of purchases at pre-existing stores per day (see Figure 3-5). We note that this series also increases steadily throughout the observation period, but that it increases faster than expected after the grocer enters the loyalty program. This indicates that consumers seem to be buying more frequently at pre-existing stores after the entry of the large grocer. We estimate that the number of daily purchases increased by 473 transactions on average. The posterior probability of a positive causal effect is 98%. The effect has a 95% Bayesian credible interval of [82, 776]. This is equivalent to a 34% increase in the observed daily purchases with a 95% Bayesian credible interval of [6%, 57%].

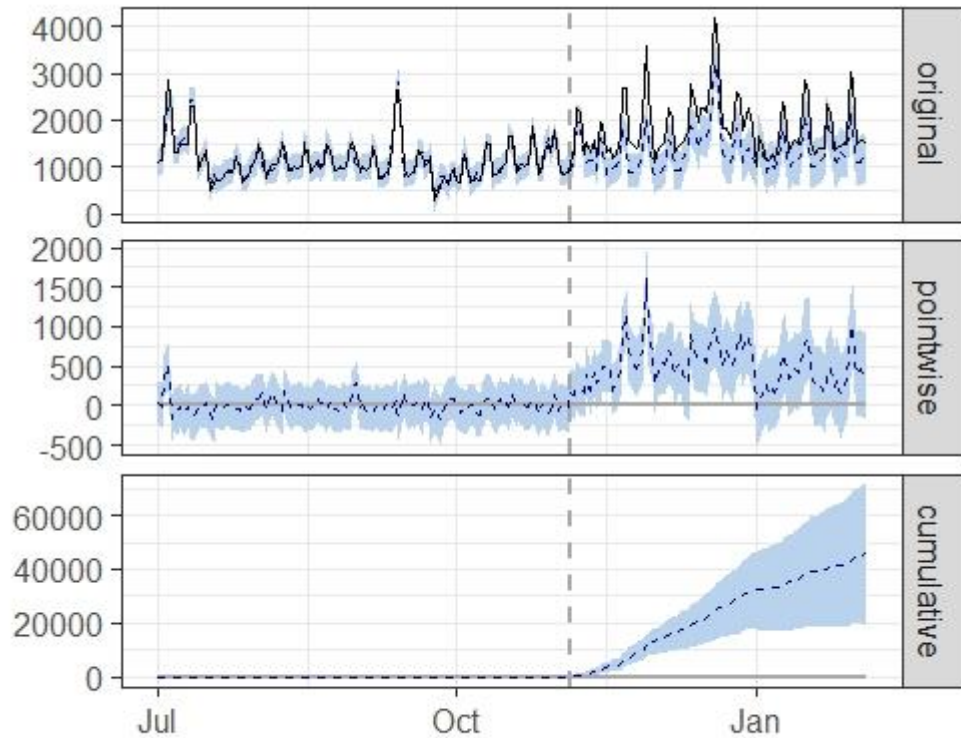


Figure 3-5 Comparing observed data to counterfactual data for total no. of transactions

Our fourth measure is the total number of points used at pre-existing stores per day (see Figure 3-6). We note that this series also increases steadily through the observation period, but that it increases faster than expected after the grocer enters the loyalty program. This indicates that consumers seem to be using points more actively at pre-existing stores after the grocer’s entry. We estimate that the number of points used at pre-existing stores increased by 0.85 points per day with a 95% Bayesian credible interval of [0.49, 1.2]. The posterior probability of a positive causal effect is 99%. This is equivalent to an 8.7% increase in the number of points used with a 95% Bayesian credible interval of [5%, 12%].

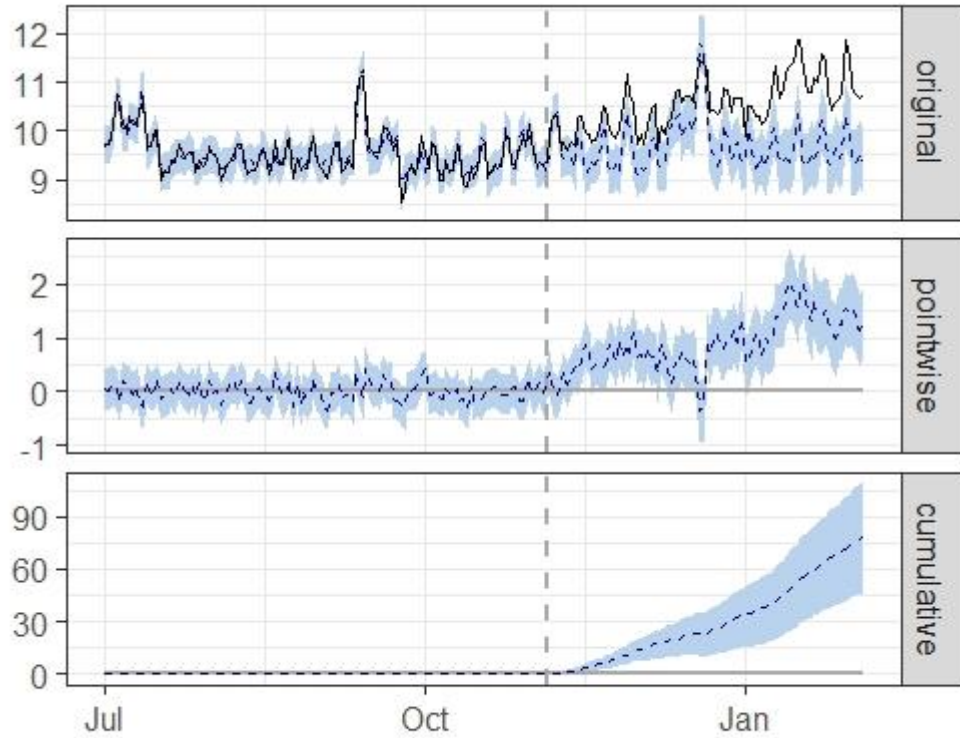


Figure 3-6 Comparing observed data to counterfactual data for points used

Our fifth and final measure is the total number of points earned at pre-existing stores per day (see Figure 3-7). We note a steadily increasing amount of points earned throughout the observation period. We estimate that the number of points earned at pre-existing stores increased by 0.83 points per day with a 95% Bayesian credible interval of [0.28, 1.4]. The posterior probability of a positive causal effect is 99%. This is equivalent to an 8.4% increase in the number of points earned with a 95% Bayesian credible interval of [2.9%, 14%].

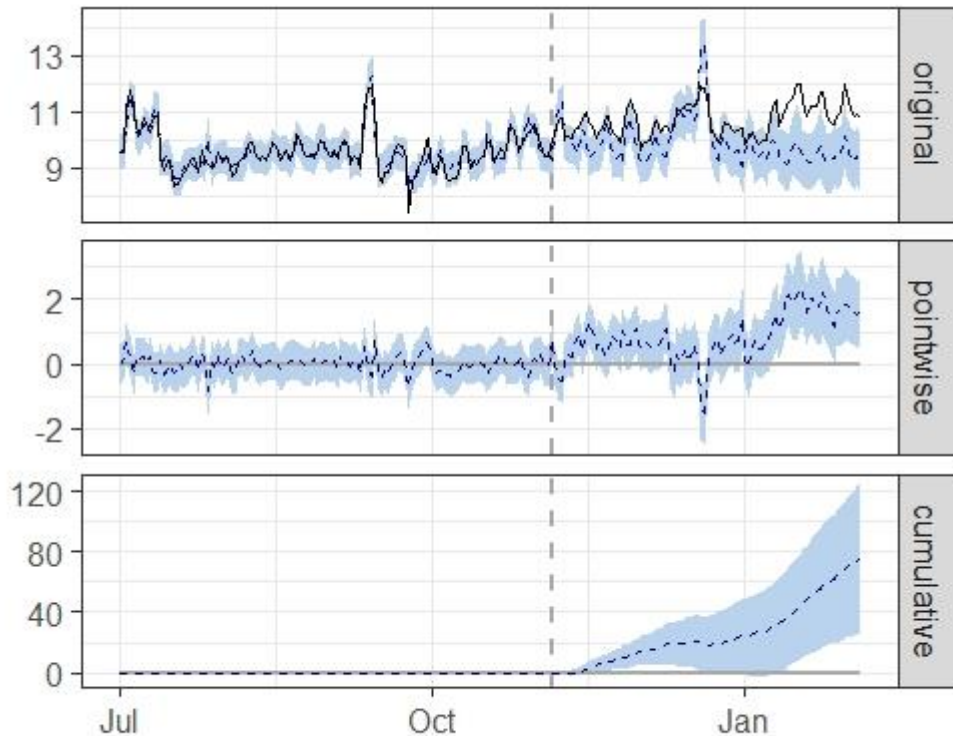


Figure 3-7 Comparing observed data to counterfactual data for points earned

These results converge with our earlier findings using an analysis of matched samples to find multiple sources of positive spillovers for coalition loyalty program merchants from new merchant entry. In particular, we see again that for the existing merchants and grocery customers, there is a statistically significant increase in aggregate sales and number of transactions. The BSTS analysis supports the previous program level studies findings of reduced average basket size and increases in both loyalty points used and earned.

3.6 Extension to other merchants

To further generalize our findings, we perform equivalent analyses for four of the next most significant merchant entries into the program, by number of customers. These entries consist of a chain of coffee shops, a group of clothing stores, a home goods vendor, and a shoe brand. These

merchant entries are significantly smaller than the grocery merchant studied in the main analysis; each retailer is about 10 percent the size of the grocer in terms of the number of stores. Summary data is shown in Table 3-11 where we again define our two customer cohorts as before.

Figure 3-8 plots the cumulative number of stores in the coalition with reference lines for each merchant entry. These merchants were chosen for the analysis because they introduce a multitude of stores into the coalition and capture at least 1,000 customers within +/- 90 days of their entry (+/- 60 in the case of the home goods retailer).

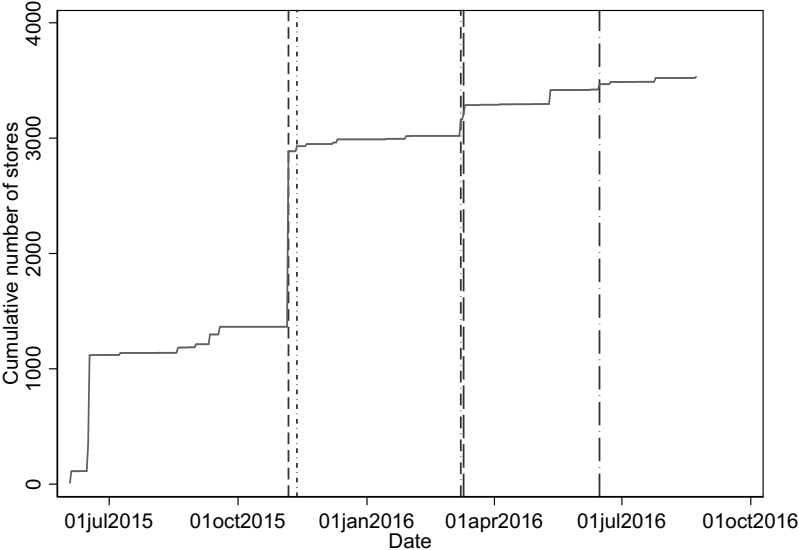


Figure 3-8 Cumulative number of stores in the coalition by merchant entry. Note: Vertical reference lines refer to the entry of large grocer, shoes/accessories retailer, café, clothing retailer, and home goods retailer from left to right, respectively.

We again use the matching methodology from section IV and arrive at the abridged results displayed in Table 3-12 where we only show the post-entry interaction with the pre-existing customers (the full results are shown in Appendix C).

Table 3-11 Merchant entry characteristics

Merchant	Entry date (MM/DD/YYYY)	Total # of stores		# of customers within +/- 90 days that	
		Within +/- 90 days of entry	Full sample	Never shop	Shop at least once
Grocer	11/6/2015	1,232	1,521	353,827	22,562
Café	3/8/2016	107	116	544,349	15,559
Home goods	6/15/2016	65	65	535,966	6,597
Clothing	3/10/2016	48	49	565,732	1,075
Shoes	11/12/2015	59	63	379,572	2,973

Note: The analysis of the entry of the home goods retailer is done using a window of +/- 60 days, since it enters within the last 90 days of our data set.

Table 3-12 Post-merchant entry interactions with relevant customers at the Program-level

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (in 000s)	(4) Points Earned	(5) Points Used
Café	5,466.18*** (644.22)	16.21** (7.85)	1,949.65*** (181.10)	270,423.28*** (21,372.74)	246,874.03*** (17,497.52)
Home goods	10,691.61*** (946.70)	9.52 (6.51)	3,011.79*** (253.32)	243,896.71*** (25,573.80)	170,136.62*** (17,521.26)
Clothing retailer	5,476.02*** (667.41)	-13.24 (12.20)	2,006.50*** (188.08)	288,205.09*** (22,151.59)	263,716.99*** (18,015.83)
Shoe retailer	1,887.07*** (421.18)	-62.59*** (17.49)	1,056.26*** (135.50)	19,617.77 (24,475.25)	-16,698.74 (14,413.53)

Robust standard errors in parentheses
 *** p<0.01, ** p<0.05, * p<0.1

We find that, similar to the grocer's entry, at the program-level these smaller merchants lead to an increase in the total number of transactions and aggregate sales at merchants that exist prior to their entry. However, we again have mixed results with regard to basket size. For the café merchant, average basket size increases, whereas for the shoe retailer it decreases. For both the home goods and clothing retailer, the effect on basket size is not statistically significant. Also, at the program-level, the café retailer, the home goods retailer, and the clothing retailers all have a positive impact on point earning and redemption behavior at other merchants, while the shoes and accessories retailer does not. At the store-level, all merchants have similar impacts across all outcome variables as the large grocer.

We find that, like the large grocer, the entries of the other merchants lead to increases in aggregate sales for all other stores, both at the program- and the store-level. This supports our finding that the introduction of new merchants into a coalition loyalty program benefits pre-existing merchants via aggregate sales. Interestingly, we find that point earning and redemption behavior is statistically greater for customers that eventually shop at the new smaller merchants entering the coalition loyalty program, and at magnitudes that are much greater than the corresponding effect from the entry of the large grocer.

3.7 Conclusion, Discussion, and Limitations

By allowing points to be earned and spent at any partner merchant, coalition loyalty programs seek to enhance outcomes for both buyers and sellers. In this paper, we measure the benefits to merchants in a coalition loyalty program when a new merchant enters the network.

Using multiple approaches, we find that shoppers for whom the entering merchant is relevant tend to spend more, and shop more frequently, at pre-existing partner merchants. Our findings show that there exist significant positive spillovers among firms in a coalition loyalty program, and that firms may do better by joining a coalition than by running a standalone loyalty program.

Whereas coalition loyalty programs in markets across the world have increased in membership and popularity over the years, little has been established about the effect of network composition on the effectiveness of such programs for each merchant's sales. The data set we use allows us to measure this effect off of merchant entries of various magnitudes into the network. Because each entering merchant is relevant to only a subset of loyalty program members, we are able to use variation in consumption habits to infer spillovers. Moreover, the panel structure of the data allows us a second empirical strategy of predicting counterfactual merchant performance had these entries not occurred. Despite the dissimilarity between these two approaches, they yield consistent findings regarding the spillovers from merchant entry on multiple outcome variables.

As with many papers in the loyalty literature, we have access to detailed transaction data for activity connected to the loyalty program, but lack data for firms prior to joining the program, as well as data for shoppers who do not participate in the program. Our empirical approach aims to address these data shortcomings; however, we concede that additional data may be helpful in identifying further types of cross-merchant spillovers.

We examine and find positive spillovers for coalition merchants from the entry of firms of various sizes and product categories. Thus, our results generalize beyond the entry of a large grocery or general merchandise firm into a coalition program. Our results may be less generalizable for coalition loyalty programs that are structured differently. Certain credit card rewards programs, for instance, share many of the same characteristics of the program we study,

e.g., the ability to earn and spend points across multiple merchants. However, unlike the program we study these programs typically have multiple firms and competitors within the same categories. Airline alliances, meanwhile, operate on the other side of the spectrum; typically allowing only one carrier in the program from each region. These programs differ from our context in several respects, and identifying spillovers from merchant entry for these cases may prove to be fertile ground for further research.

Many open questions exist pertaining to coalition loyalty programs. The composition of product categories in a program may, for instance, have significant implications on its performance. Likewise, the “exchange rates” of points between different merchants are likely important variables that determine the relative performance of each merchant in the network. These are exciting areas for future research.

References

- Abdulkadiroglu A, Sonmez T (1998). Random Serial Dictatorship and the Core from Random Endowments in House Allocation Problems. *Econometrica*. 66 (3): 689.
- Agresti A (2015) Foundations of linear and generalized linear models. Wiley series in probability and statistics.
- Anderson BA, Laurent PA, Yantis S (2011) Value-driven attentional capture. *Proceedings of the National Academy of Sciences*. 108(25):10367–10371.
- Balseiro S, Feldman J, Mirrokni V, and Muthukrishnan S (2014) Yield Optimization of Display Advertising with Ad. *Management Science*. 60(12):2886-2907.
- Benway JP, Lane DM (1998) Banner blindness: Web searchers often miss “obvious” links. *Internetworking ITG Newsletter*. 1(3):1–22.
- Bertrand M, Duflo E, Mullainathan S (2004) How much should we trust differences-in-differences estimates? *The Quarterly Journal of Economics*. Volume 119, Issue 1, 249–275.
- Blaustein A (2017) Programmatic Buying 101: The Cost Efficiencies. from <https://katana.media/blog/programmatic-buying-101-value-cost-efficiencies/>
- Brodersen KH, Gallusser F, Koehler J, Remy N, Scott SL (2015) Inferring Causal Impact using Bayesian Structural Time Series. *The Annals of Applied Statistics*. Vol. 9, No. 1, 247–274.
- Brebion A (2018, February 4) Above the Fold vs. Below the Fold: Does it Still Matter in 2019. from <https://www.abtasty.com/blog/above-the-fold/>
- Casella A (2003) Storable Votes. *Games and Economic Behavior*. 51:391–419.
- Chatterjee P, Hoffman D, Novak TP (2003) Modeling the Clickstream: Implications for Web-Based Advertising Efforts. *Marketing Science*. 22(4):520–541.
- Cho CH, Cheon HJ (2004) Why do people avoid advertising on the internet? *Journal of Advertising*. 33(4):89–97.
- Cochran WG (1968) The effectiveness of adjustment by subclassification in removing bias in observational studies. *Biometrics*. 24(2):295-313.
- Cochran WG, Rubin DB (1973), Controlling Bias in Observational Studies: A Review. *Sankhyā: The Indian Journal of Statistics, Series A*. 35(4):417-446.
- Cole R, Gkatzelis V, and Goel G (2012) Mechanism Design for Fair Division. arXiv:1212.1522 <http://arxiv.org/abs/1212.1522>
- Cooper R, Shallice T (2000) Contention Scheduling and the Control of Routine Activities. *Cognitive Neuropsychology*. 17(4):297–338.
- Cramton P, Gibbons R, and Klemperer P (1987) Dissolving a partnership efficiently. *Econometrica*. 55, 615–632.

- Danaher PJ, Mullarkey GW (2003) Factors Affecting Online Advertising Recall: A Study of Students. *Journal of Advertising Research*. 43(3):252–267.
- Danaher PJ, Rust RT (1996) Determining the optimal return on investment for an advertising campaign. *European J. Oper. Res.* 95: 511–521.
- Danaher PJ, Sajtos L, Danaher TS (2017) Dynamic Management of Rewards in Coalition Loyalty Programs. *SSRN Working Paper*.
- Dehejia RH, Wahba S (1999) Causal Effects in Nonexperimental Studies: Reevaluating the Evaluation of Training Programs. *Journal of the American Statistical Association*. 94(448):1053-1062.
- Devanur NR, Hayes TP (2009) The adwords problem: Online keyword matching with budgeted bidders under random permutations. *ACM EC*. 2009.
- Dimitrov NC, Plaxton CG (2008) Competitive weighted matching in transversal matroids. *ICALP* 1:397-408.
- Dorotic M, Fok D, Verhoe PC, Bijmolt THA (2011) Do vendors benefit from promotions in a multi-vendor loyalty program? *Marketing Letters* 22: 341-356.
- Dreze X, Hussherr FX (2003) Internet advertising: Is anybody watching? *Journal of Interactive Marketing*. 17(4):8–23.
- Dreze X, Nunes JC (2004) Using Combined-Currency Prices to Lower Customers' Perceived Cost. *Journal of Marketing Research*. Vol. XLI, pp. 59-72.
- Feldman J, Henzinger M, Korula N, Mirrokni VS, Stein C (2018) Online Stochastic Packing Applied to Display Ad Allocation. <https://arxiv.org/abs/1001.5076>
- Fisher L (2018) US Programmatic Ad Spending Forecast Update 2018. <https://www.emarketer.com/content/us-programmatic-ad-spending-forecast-update-2018>
- Fleishman EA, Quaintance MK (1984) *Taxonomies of Human Performance: The Description of Human Tasks*. (Academic Press Inc.).
- Forster S, Lavie N (2009) Harnessing the wandering mind: the role of perceptual load. *Cognition*. 111(3):345-55.
- Gelman A, Carlin JB, Stern HS, Rubin DB (2002) *Bayesian Data Analysis*, 2nd ed. Boca Raton, FL: Chapman and Hall/CRC.
- George EI, McCulloch RE (1997) Approaches for Bayesian variable selection. *Statist. Sinica*. 7: 339–374.
- Gollapudi S, Panigrahi D (2014) Fair Allocation in Online Markets. *ACM CIKM*, 2014.
- Ha, A (2019) eMarketer predicts digital ads will overtake traditional spending in 2019. <https://techcrunch.com/2019/02/20/emarketer-digital-ad-forecast/>
- Hansen CB (2007) Generalized least squares inference in panel and multilevel models with serial correlation and fixed effects. *Journal of Econometrics*. 140: 670–694.

- Hastie T, Tibshirani R, Friedman J (2009) *The Elements of Statistical Learning Data Mining, Inference, and Prediction*. (Second Edition. Springer).
- Heine C (2014, December 24) Adweek, Agency and Brand Leaders Weigh in on Digital's Mounting Ad Viewability Issues. Retrieved from <https://www.adweek.com/digital/agency-and-brand-leaders-weigh-digital-ad-viewability-issues-162082/>
- Hoban PR, Bucklin RE (2015) Effects of Internet Display Advertising in the Purchase Funnel: Model-Based Insights from a Randomized Field Experiment. *Journal of Marketing Research*. 52(3):375–393.
- Hobbs T (2016) Nectar plans to become the UK's biggest digital loyalty brand. *Marketing Week*. Accessed online on 11 Dec 2018 <<https://www.marketingweek.com/2016/03/09/how-nectar-plans-to-become-the-biggest-digital-loyalty-brand-in-the-uk/>>.
- Huberman BA, Pirolli PLT, Pitkow JE, Lukose RM (1998) Strong regularities in World Wide Web surfing. *Science*. 280(5360):95–97.
- Ilfeld JS, Winer RS (2002) Generating Web Site Traffic: An Empirical Analysis of Web Site Visitation Behavior. *Journal of Advertising Research*. 42(5):49-61.
- Jackson MO, Sonnenschein HF (2007) Overcoming Incentive Constraints by Linking Decisions. *Econometrica*. Vol. 75 (1).
- Janiszewski C, Kuo A, Tavassoli N (2013) The Influence of Selective Attention and Inattention to Products on Subsequent Choice. *Journal of Consumer Research*. 39(6): 1258-1274.
- Jiang Z, Nevskaya Y, Thomadsen R (2017) Can Non-Tiered Customer Loyalty Programs be Profitable? *SSRN Working Paper*.
- Kalman RE (1960) A New Approach to Linear Filtering and Prediction Problems. *Transactions of the ASME--Journal of Basic Engineering*. 82(D):35-45.
- Kalyanasundaram B, Pruhs KR (2000) An optimal deterministic algorithm for online b-matching. *Theoretical Computer Science*. 233(1–2):319-325.
- Kennedy J, Eberhart R (1995) Particle swarm optimization. *Proceedings of IEEE International Conference on Neural Networks*. pages 1942-1948.
- Kim TH, Maruta I, Sugie T (2008) Robust PID controller tuning based on the constrained particle swarm optimization. *Automatica*. 44:1104-1110.
- Kopalle PK, Sun Y, Neslin S, Sun B, Swaminathan V (2012) The Joint Sales Impact of Frequency Reward and Customer Tier Components of Loyalty Programs. *Marketing Science*. 31(2): 216-235.
- Koshevoy GA, Mosler K (1997) Multivariate Gini Indices. *Journal of Multivariate Analysis*. 60:252-276.
- Kwong JYY, Soman D, Ho CKY (2011) The role of computational ease on the decision to spend loyalty program points. *Journal of Consumers Psychology*. 21: 146-156.

- Lamy D, Leber A, Egeth HE (2004) Effects of Task Relevance and Stimulus- Driven Saliency in Feature-Search Mode. *Journal of Experimental Psychology: Human Perception and Performance*. 30(6):1019–1031.
- Lavie N (1995) Perceptual load as a necessary condition for selective attention. *Journal of Experimental Psychology: Human Perception and Performance*. 21(3):451–468.
- Lavie N, Tsai Y (1994) Perceptual load as a major determinant of the locus of selection in visual attention. *Perception & Psychophysics*. 56(2):183–197.
- Lewis M (2004) The Influence of Loyalty Programs and Short-Term Promotions on Customer Retention. *Journal of Marketing Research*, Vol. XLI, 281-292.
- Li H, Kannan PK (2014) Attributing Conversions in a Multichannel Online Marketing Environment: An Empirical Model and a Field Experiment. *Journal of Marketing Research*. 51(1):40–56.
- Li H, Kannan PK, Viswanathan S, Pani P (2016) Attribution Strategies and Return on Keyword Investment in Paid Search Advertising. *Marketing Science*. 35(6):831-848.
- Liu Y (2007) The Long-Term Impact of Loyalty Programs on Consumer Purchase Behavior and Loyalty. *Journal of Marketing*. 71: 19-35.
- Manchanda P, Dubé J, Goh KY, Chintagunta PK (2006) The Effect of Banner Advertising on Internet Purchasing. *Journal of Marketing Research*. 43(1):98–108.
- Manshadi VH, Gharan SO, Saberi A (2010) Online Stochastic Matching: Online Actions Based on Offline Statistics. *ACM SIAM*, 2010.
- McAfee P (1992) Amicable Divorce: Dissolving a Partnership with Simple Mechanisms. *Journal of Economic Theory*. 56.
- Mehta A (2012) Online Matching and Ad Allocation. *Foundations and Trends in Theoretical Computer Science*. 8(4): 265-368.
- Moe, WA (2006) A field experiment to assess the interruption effect of pop-up promotions. *Journal of Interactive Marketing*. 20(1):34–44.
- Moore RS, Stammerjohan CA, Coulter RA (2005) Banner advertiser-web site congruity and color effects on attention and attitudes. *Journal of Advertising*. 34(2):71–84.
- Nectar (2018) <<https://www.nectar.com>>.
- Norman DA, Shallice T (1986) Attention to action. In *Consciousness and self-regulation* (pp. 1-18). Springer, Boston, MA.
- Norman DA, Shallice T (1980) Attention to Action: Willed and Automatic Control of Behavior Technical Report No. 8006.
- Orhun AY, Guo T (2018) Reaching for Gold: Frequent-Flyer Status Incentives and Moral Hazard. Available at SSRN: <https://ssrn.com/abstract=3289321>
- Pan F, Li S, Ao X, Tang P, He Q (2019) “Warm Up Cold-start Advertisements: Improving CTR Predictions via Learning to Learn ID Embeddings” *ACM SIGIR*.

- Payback (2018) <<https://www.payback.net/en/>>.
- Pearson B (2018) 4 Ways The Demise Of Plenti Will Go On To Reward Shoppers. *Forbes*. Accessed online on 12 Dec 2018 < <https://www.forbes.com/sites/bryanpearson/2018/05/07/4-ways-the-demise-of-plenti-will-go-on-to-reward-shoppers/#4465195c33d6>>.
- Perlich C, Dalessandro B, Hook R, Stitelman O, Raeder T, Provost F (2012) Bid optimizing and inventory scoring in targeted online advertising. *ACM KDD*, pages 804-812.
- Rossi F (2017) Low Price or Higher Reward? Measuring the Effect of Consumers' Preferences on Reward Programs. *Management Science*. 64(9): 3971-4470.
- Rubin DB (1974) Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*. 66(5):688-701.
- Rutz OJ, Bucklin RE (2012) Does banner advertising affect browsing for brands? clickstream choice model says yes, for some. *Quantitative Marketing and Economics*. 10(2):231–257.
- Santos A, Anta AF, Cuesta JA, Fernández LL (2016) Fair linking mechanisms for resource allocation with correlated player types. *Computing*. 98(8):777–801.
- Scott SL (2010) A modern Bayesian look at the multi-armed bandit. *Applied Stochastic Models in Business and Industry*. 26(6):639–658.
- Scott SL, Varian HR (2015) Bayesian Variable Selection for Nowcasting Economic Time Series. Chapter in NBER book *Economic Analysis of the Digital Economy*, Avi Goldfarb, Shane M. Greenstein, and Catherine E. Tucker, editors (p. 119 - 135).
- Shapiro B (2018) Positive Spillovers and Free Riding in Advertising of Prescription Pharmaceuticals: The Case of Antidepressants. *Journal of Political Economy*. 126(1).
- Shoulberg W (2018) Plenti Customer Loyalty Program to End in July: What Went Wrong? *Forbes*. Accessed online on 11 Dec 2018 <<https://www.forbes.com/sites/warrenshoulberg/2018/05/01/no-good-and-plenti-as-the-customer-loyalty-program-winds-down/#678055ef156d>>.
- Sports Loyalty International, Inc. (2018) Coalition Loyalty: A Model with Sustainable Advantages for Retailers. *SLI White Paper*. Accessed online on 12 Dec 2018 <<http://www.sli21.com/?whitepaper=coalition-loyalty-a-model-with-sustainable-advantages-for-retailers>>.
- Stourm V, Bradlow ET, Fader PS (2015) Stockpiling Points in Linear Loyalty Programs. *Journal of Marketing Research*. Vol. LII, pp. 253-267.
- Stourm V, Bradlow ET, and Fader PS (2018) Market Positioning using Cross-Reward Effects in a Coalition Loyalty Program. *SSRN Working Paper*.
- Strong EC (1977) The Spacing and Timing of Advertising. *Journal of Advertising Research*. 17(6):25-31.
- Sung SW (2009) *Process Identification and PID Control*. Wiley-IEEE Press.

- Tavassoli N, Shultz CK, Fitzsimons GJ (1995) Program involvement: Are more moderate levels best for ad memory and attitude toward the ad? *Journal of Advertising Research*. 35(5):61-73.
- Trusov M, Ma L, Jamal Z (2016) Crumbs of the Cookie: User Profiling in Customer-Base Analysis and Behavioral Targeting. *Marketing Science*. 35(3):405-426.
- Van de Heere HJ, Bijmolt THA (2005) Decomposing the Promotional Revenue Bump for Loyalty Program Members Versus Nonmembers. *Journal of Marketing Research*, Vol. LII, pp. 443-457.
- Wei L, Xiao J (2015) Are points like money? An empirical investigation of reward promotion effectiveness for multicategory retailers. *Market Letters*. 26: 99-114.
- West M, Harrison J (1997) Bayesian forecasting and dynamic models. Springer series in statistics. 2nd edition.
- Work S, Hayes M (2018) Buy Banner Ads. Retrieved from <https://www.shopify.com/guides/make-your-first-ecommerce-sale/banner-ads>
- Yi G, Jeong HM, Choi W, Jang S, Lee H, Kim BJ (2014) Human dynamics of spending: Longitudinal study of a coalition loyalty program. *Physica A* 310: 391-398.
- Zellner A (1986) On assessing prior distributions and Bayesian regression analysis with g-prior distributions. In Bayesian Inference and Decision Techniques: Essays in Honor of Bruno de Finetti (P. K. Goel and A. Zellner, eds.). Stud. Bayesian Econometrics Statist. 6 233-243. North-Holland, Amsterdam.
- Zhang J, Breugelmans E (2012) The Impact of an Item-Based Loyalty Program on Consumer Purchase Behavior. *Journal of Marketing Research*, Vol. XLIX, pp. 50-65.
- Zhang W, Rong Y, Wang J, Zhu T, Wang X (2016) Feedback Control of Real-Time Display Advertising. *ACM WSDM*, 2016.
- Zhang J, Wedel M (2009) The Effectiveness of Customized Promotions in Online and Offline Stores. *Journal of Marketing Research*, Vol. XLVI, pp. 190-206.
- Ziegler JG, Nichols NB (1993) Optimum settings for automatic controllers. *Journal Of Dynamic Systems Measurement And Control-Transactions Of The Asme*. 115(2B):220-222.

Appendix

A Online Task Progression and Display Ad Engagement

A.1 Propensity Score Matching Details

When performing matching, we make use of an interval around the score being matched to prevent us from matching observations that are not similar enough. This is standard practice and recommended by Cochran and Rubin (1973).

The matching algorithm in dataset 1 partitions the dataset into three pieces based on whether the observations were first, middle or last steps. We seek to match first to middle and last to middle by sampling without replacement. We start by taking all the first step observations and randomly picking one and match it to a middle step observation that has the nearest linearized first step propensity score that is at most 0.5 away in absolute value. Our linearized propensity scores tend to range from -1 to 2. Cochran (1968) has shown that 90% of the variation of most distributions can be captured by just five equal partitions. We therefore use intervals of 0.5 to ensure the equivalent of at least five such partitions for the propensity score distributions. These two matched observations are then removed from their partitions and placed into a first step matched dataset. The procedure is then repeated until there are no more first or middle steps remaining to create more matched pairs. This entire matching process is repeated with the last and middle steps thus yielding our two balanced datasets on which we can estimate the logit model previously defined. An identical algorithm is implemented in dataset 2 where the only change is that we have more treatment groups.

A.2 SEC DSP Cost Estimates

Below are select operating costs from three DSPs (MaxPoint, Rocket Fuel and Criteo) who have public record of spending in 2016 (the year we collected our data). The average media cost percentage is 46% of revenue and the average computation running costs (reflected in R&D expenditure) is 11% of revenue.

Table A-1 DSP Cost of Operations

2016 Operations Filings (in thousands)			
	MXPT	FUEL	CRTO
Revenue	\$149,109	\$456,263	\$1,799,146
Media Cost	\$51,120	\$204,168	\$1,068,911
R&D	\$26,576	\$35,354	\$123,649
Media Cost percentage	34%	45%	59%
R&D percentage	18%	8%	7%

B Real-time Digital Ad Allocation: A Fair Streaming Allocation

Mechanism

B.1 PID results without PSO

In this section, we add a massive supply shock to the data set to illustrate the need for the particle swarm optimization. From minutes 400 through 600, we increase the ad supply by an order of magnitude (10 times). This not highly unusual and certainly something that any real-world solution would need to be able to handle. We then compare the algorithm as presented in the paper with one where the particle swarm optimization component is removed. For this second model, we set the PID parameters to the stable set found in the first simulation of this method.

In Figure B-1, we show the output of the algorithm where it rapidly adjusts to the shock, raises the threshold, and still provides sufficient ads at high performance values. In Figure B-2, we disable the particle swarm optimization and note that the threshold takes substantially longer to adjust to changes in the supply environment yielding extended time periods with far too many or too few ad spaces being won.

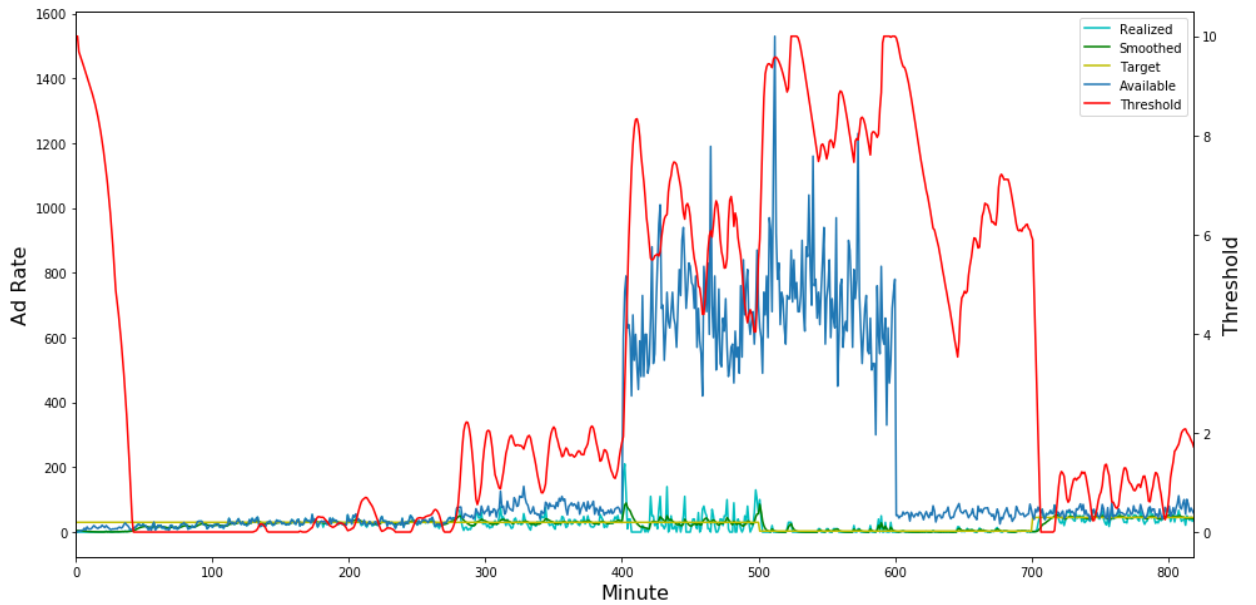


Figure B-1 Agent Level Results of Proposed Method

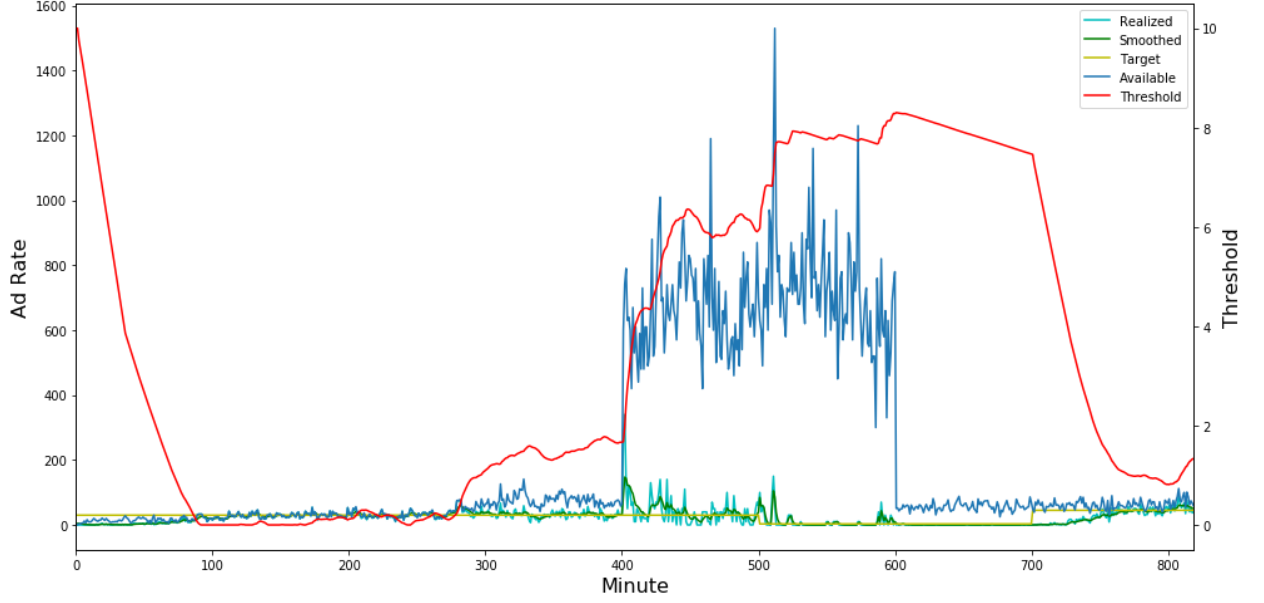


Figure B-2 Agent Level Results of Proposed Method without Particle Swarm Optimization

B.2 Market level DH formulation

The adwords solution proposed by Devanur and Hayes (2009) is a close but not perfect algorithm for our particular problem and we have adjusted it to be a useful benchmark. We follow their implementation directly where we select the first ε (in our case 10%) of the observations as a training set. Per their design, these training items are assigned randomly to advertisers. Then we learn a set of weights $\alpha^* := \operatorname{argmin}_{\alpha} \{D(\alpha, S)\}$ for each advertiser i , where S is the set the observed items j with preference score r_{ij} . And $D(\alpha, S) := \sum_i \alpha_i \varepsilon B_i + \sum_j \max_i r_{ij} (1 - \alpha_i)$. The algorithm allocates item j to advertiser i that maximizes $r_{ij} (1 - \alpha_i)$.

A direct conversion of the Devanur and Hayes (2009) method would create a new budget according to anticipated future items. This would be done by summing the observed preference scores for each advertiser during the training period and then dividing by the number of advertisers and

then adjust for the remaining $(1 - \epsilon)$ items. That is, we set $B_i = \frac{\sum_j r_{ij}}{\# \text{ of agents}} \times \frac{1-\epsilon}{\epsilon}$. However, since the budget has no practical use we do not enforce it. Intuitively, adding such a budget adds fictional constraints which only decrease efficiency. This can easily be confirmed via simulation.

It is worth mentioning that the Devanur and Hayes (2009) implementation does not approach greedy under any non-trivial circumstances. While the allocation problem is similar to ours, the search problem is not. The primary objective of a search engine is to exhaust all advertiser budgets. The Devanur and Hayes (2009) algorithm will always exhaust advertiser budgets, which often requires allocating items to a lower scoring advertiser for the sake of exhausting that advertiser’s budget. The goal of their method is to achieve that highest efficiency possible given that it *must exhaust all budgets*. Our implementation is centered on advertiser quota and not advertiser budget and so it does not directly have this constraint. We are able to achieve higher levels of efficiency because we can simultaneously select more low scoring items and use a smoother fairness tradeoff between advertisers competing for higher scoring items.

C Cross-Merchant Spillovers in Coalition Loyalty Programs

C.1 Variable Definitions

Table C-1 Variable definitions for difference-in-differences analysis

Name	Definition
Total no. of transactions	The total number of transactions that occur on a given day
Average basket size	A simple average of all basket sizes on a given day, net of all discounts applied (i.e. points used plus applied multipliers)
Aggregate sales	The sum all basket sizes on a given day, net of all discounts applied (i.e. points used plus applied multipliers)

Point earned	The sum of all points earned from campaign-related transactions, regardless of campaign type (i.e. campaigns that have a fixed amount of points that can be earned versus a campaign where points are earned as a percent of basket size)
Points used	The sum of all points used for campaign-related transactions, not including additional discounts from applied multipliers
Post-merchant entry	An indicator for whether the new merchant has entered the program based on that merchant's first transaction in the data
Treated customers	An indicator for whether the observation applies to all customers that eventually shop at the new merchant
Post-merchant entry * Treated customers	An interaction term of Post-merchant entry and Treated customers, i.e. the DD estimator

Table C-2 Variable definitions for propensity score matching analysis

Name	Definition
Aggregate sales	The sum of all basket sizes by customer over the 90 day period prior to the large merchant's entry. As before, aggregate sales is calculated net of all points used and discounts applied.
Average basket size	A simple average of all basket sizes by customer over the 90 day period prior to the large merchant's entry, net of all points used and discounts applied.
Number of stores visited	The number of unique stores visited by each customer over the 90 day period prior to the large merchant's entry
Number of days making a transaction	The number of unique days on which an individual makes at least one transaction over the 90 day period prior to the large merchant's entry.

Table C-3 Variable definitions for Bayesian structural time series

Variable	Shortened variable name for cross-references	Notes
----------	--	-------

Non-grocery customer points earned	ce_control	sum of non-grocery store points earned by customers that day that never make a purchase at a grocery store
Non-grocery customer points used	cu_control	sum of non-grocery store points used by customers that day that never make a purchase at a grocery store
total purchase amount control	fp_control	total sales of all transactions at non-grocery stores (including points used) by day for non-grocery store shoppers
Non-Grocery customer purchases control	n_purch_control	the total number of non-grocery store purchases made by customers that day that never make a purchase at a grocery store
Shriver and Bolinger price index	price_index_sb	Shriver and Bolinger price index

C.2 Matching Analysis

In this appendix we repeat the analysis of section IV using the entire unmatched sample instead of the propensity score-matched sample. We find that these results are generally consistent with our main analysis.

Table C-4 Main results (aggregate)

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (in 000s)	(4) Points earned	(5) Points used
Post-Merchant entry	-3,667.52*** (349.88)	60.24*** (10.27)	-1,112.04*** (107.58)	-27,923.41 (22,567.46)	2,601.66 (13,179.88)
Treated customers	-7,191.37*** (348.88)	-22.06** (9.71)	-2,797.04*** (99.05)	-126360.44*** (21,584.05)	-107962.22*** (11,658.37)
Post * Treated	3,683.41*** (370.00)	-60.46*** (13.97)	1,123.57*** (116.85)	34,253.22 (22,907.08)	5,485.05 (13,523.62)
Constant	7,848.75*** (338.18)	402.49*** (6.77)	3,044.82*** (93.38)	137,573.53*** (21,382.25)	116,711.77*** (11,441.37)
Observations	362	362	362	362	362
R-squared	0.76	0.22	0.83	0.25	0.44
Day-of-week FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

Table C-5 Main results (store level)

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (in 000s)	(4) Points earned	(5) Points used
Post-Merchant entry	-6.15*** (0.10)	-105.90*** (2.64)	-1.87*** (0.05)	-50.94*** (8.71)	17.20*** (6.30)
Treated customers	-11.25*** (0.11)	-246.40*** (2.57)	-4.31*** (0.04)	-342.82*** (9.74)	-312.17*** (5.98)
Post * Treated	6.48*** (0.12)	94.33*** (3.20)	2.05*** (0.06)	74.48*** (9.91)	20.00*** (7.69)
Constant	13.98*** (0.16)	391.65*** (3.23)	5.19*** (0.06)	426.32*** (14.94)	359.69*** (9.05)
Observations	238,819	238,819	238,819	115,797	115,797
R-squared	0.48	0.25	0.48	0.09	0.16
Day-of-week FE	YES	YES	YES	YES	YES
Store FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

Table C-6 Heterogeneity (aggregate)

VARIABLES	(1) Total # of transactions	(2) Average basket size	(3) Aggregate sales (000s)	(4) Points Earned	(5) Points Used
Post-Merchant entry	-3,669.53*** (377.02)	60.28*** (10.29)	-1,112.48*** (120.12)	-27,951.88 (22,888.31)	2,556.69 (13,577.16)
Low-value Merchant customers	-7,648.66*** (362.80)	-46.48*** (11.84)	-2,974.32*** (103.82)	-133729.79*** (21,659.14)	-114028.49*** (11,736.56)
Mid-value Merchant customers	-3,814.48*** (181.44)	-14.11** (5.49)	-1,481.41*** (51.92)	-67,003.46*** (10,829.50)	-56,871.66*** (5,868.24)
High-value Merchant customers	-2,538.43*** (120.97)	0.95 (3.67)	-983.47*** (34.62)	-44,609.34*** (7,219.41)	-37,901.46*** (3,912.50)
Post- Merchant entry * Low-value Merchant Customers	3,666.05*** (382.21)	-75.40*** (16.32)	1,110.06*** (122.45)	29,627.73 (22,977.48)	-57.09 (13,665.94)
Post- Merchant entry * Mid-value Merchant Customers	1,833.60*** (191.19)	-26.06*** (7.95)	557.22*** (61.25)	14,709.45 (11,489.62)	-265.10 (6,833.92)
Post- Merchant entry * High-value Merchant Customers	1,231.08*** (127.45)	-18.42*** (5.19)	374.96*** (40.82)	10,389.09 (7,659.59)	349.63 (4,556.94)
Constant	7,849.77*** (360.04)	402.47*** (6.81)	3,045.04*** (102.41)	137,587.84*** (21,607.51)	116,734.38*** (11,679.45)
Observations	724	724	724	724	724
R-squared	0.81	0.18	0.86	0.32	0.54
Day-of-week FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

Table C-7 Heterogeneity (store level)

VARIABLES	(1) Total # of	(2) Average basket	(3) Aggregate	(4) Points Earned	(5) Points Used
-----------	-------------------	-----------------------	------------------	----------------------	--------------------

	transactions	size	sales (000s)		
Post- Merchant entry	-4.91*** (0.09)	-115.10*** (2.45)	-1.58*** (0.05)	-50.27*** (4.29)	-15.32*** (3.12)
Low-value Merchant customers	-9.72*** (0.09)	-290.59*** (2.33)	-3.79*** (0.04)	-172.11*** (3.70)	-148.76*** (2.14)
Mid-value Merchant customers	-4.86*** (0.04)	-144.38*** (1.15)	-1.89*** (0.02)	-86.61*** (1.86)	-74.50*** (1.07)
High-value Merchant customers	-3.26*** (0.03)	-95.82*** (0.81)	-1.26*** (0.01)	-58.09*** (1.24)	-50.05*** (0.72)
Post-Merchant entry * Low-value Merchant Customers	5.04*** (0.10)	108.18*** (2.80)	1.65*** (0.05)	55.38*** (4.32)	22.07*** (3.14)
Post-Merchant entry * Mid-value Merchant Customers	2.53*** (0.05)	55.21*** (1.39)	0.83*** (0.02)	27.81*** (2.16)	10.96*** (1.56)
Post-Merchant entry * High-value Merchant Customers	1.70*** (0.03)	37.94*** (0.93)	0.56*** (0.01)	19.32*** (1.44)	8.0691*** (1.04)
Constant	10.52*** (0.11)	349.96*** (2.38)	4.02*** (0.04)	196.01*** (5.02)	166.05*** (2.98)
Observations	588,667	588,667	588,667	588,667	588,667
R-squared	0.28	0.19	0.26	0.06	0.10
Day-of-week FE	YES	YES	YES	YES	YES
Store FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

C.3 Other merchant entries

In this appendix, we consider additional merchant entries to examine whether positive spillovers experienced from the large grocer's entry extend to smaller merchants entering the coalition. The following tables replicate the matching analysis from section IV for a home goods store, café, clothing retailer and a shoe retailer that each joined the coalition loyalty program.

Table C-8 Café retailer entry, program-level analysis

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (in 000s)	(4) Points Earned	(5) Points Used
Post-merchant entry	-5,471.02*** (590.15)	-49.12*** (5.71)	-1,974.71*** (165.23)	-278161.70*** (20,389.60)	-253684.02*** (16,464.39)
Café customers	-20,432.58*** (584.23)	-11.70 (7.23)	-4,628.37*** (167.00)	-374851.91*** (20,363.20)	-358237.48*** (16,602.19)
Post-merchant entry * Café customers	5,466.18*** (644.22)	16.21** (7.85)	1,949.65*** (181.10)	270,423.28*** (21,372.74)	246,874.03*** (17,497.52)
Constant	21,260.39*** (555.31)	226.68*** (5.46)	4,805.03*** (158.74)	389,099.15*** (19,886.20)	371,784.90*** (16,098.17)

Observations	362	362	362	362	362
R-squared	0.90	0.32	0.85	0.72	0.78
Day-of-week FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

Table C-9 Café retailer entry, store-level analysis

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (in 000s)	(4) Points Earned	(5) Points Used
Post-merchant entry	-2.71*** (0.06)	-25.36*** (0.84)	-0.94*** (0.02)	-170.90*** (3.77)	-157.46*** (3.39)
Café customers	-11.14*** (0.06)	-116.12*** (0.81)	-2.49*** (0.02)	-404.29*** (7.96)	-403.55*** (7.34)
Post-merchant entry * Café customers	2.78*** (0.07)	20.65*** (1.08)	0.94*** (0.03)	192.62*** (6.77)	179.10*** (6.45)
Constant	12.24*** (0.08)	159.17*** (0.89)	2.77*** (0.03)	336.23*** (6.53)	326.05*** (5.91)
Observations	668,592	668,592	668,592	336,204	336,204
R-squared	0.40	0.29	0.42	0.21	0.24
Day-of-week FE	YES	YES	YES	YES	YES
Store FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

Table C-10 Home goods retailer entry, program-level analysis

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (in 000s)	(4) Points Earned	(5) Points Used
Post-merchant entry	-10,731.84*** (884.85)	-53.54*** (4.13)	-3,041.12*** (237.87)	-246755.22*** (24,870.00)	-170771.07*** (16,823.24)
Home goods customers	-26,974.80*** (815.56)	-18.91*** (3.81)	-5,455.75*** (214.63)	-332561.33*** (23,591.37)	-278057.65*** (15,038.03)
Post-merchant entry * Home goods customers	10,691.61*** (946.70)	9.52 (6.51)	3,011.79*** (253.32)	243,896.71*** (25,573.80)	170,136.62*** (17,521.26)
Constant	27,431.04*** (779.07)	199.03*** (2.30)	5,540.65*** (205.54)	337,649.53*** (23,211.63)	282,312.34*** (14,628.42)
Observations	242	242	242	242	242
R-squared	0.91	0.54	0.85	0.67	0.75
Day-of-week FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

Table C-11 Home goods retailer entry, store-level analysis

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (in 000s)	(4) Points Earned	(5) Points Used
-----------	-------------------------------------	-------------------------------	-------------------------------------	----------------------	--------------------

Post-merchant entry	-4.77*** (0.11)	-38.01*** (0.84)	-1.29*** (0.02)	-121.98*** (2.54)	-83.22*** (1.91)
Home goods customers	-13.44*** (0.12)	-123.29*** (0.81)	-2.68*** (0.03)	-349.11*** (6.51)	-309.47*** (5.10)
Post-merchant entry * Home goods customers	4.85*** (0.12)	33.78*** (1.04)	1.30*** (0.03)	156.92*** (6.26)	121.51*** (5.47)
Constant	13.90*** (0.12)	146.46*** (0.86)	2.78*** (0.03)	233.42*** (4.45)	199.33*** (3.11)
Observations	493,776	493,776	493,776	244,511	244,511
R-squared	0.25	0.30	0.38	0.22	0.30
Day-of-week FE	YES	YES	YES	YES	YES
Store FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

Table C-12 Clothing retailer entry, program-level analysis

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (in 000s)	(4) Points Earned	(5) Points Used
Post-merchant entry	-5,464.61*** (609.78)	-48.52*** (5.36)	-2,004.97*** (170.47)	-288263.29*** (21,062.92)	-263555.69*** (16,862.33)
Clothing retailer customers	-22,292.67*** (603.03)	19.15* (10.57)	-4,968.57*** (172.66)	-403458.73*** (21,113.03)	-386277.69*** (17,115.77)
Post-merchant entry * Clothing retailer customers	5,476.02*** (667.41)	-13.24 (12.20)	2,006.50*** (188.08)	288,205.09*** (22,151.59)	263,716.99*** (18,015.83)
Constant	22,328.0211*** (571.44)	222.94*** (5.07)	4,976.62*** (163.18)	404,248.51*** (20,543.24)	386,969.35*** (16,511.07)
Observations	362	362	362	362	362
R-squared	0.91	0.23	0.85	0.73	0.78
Day-of-week FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

Table C-13 Clothing retailer entry, store-level analysis

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (in 000s)	(4) Points Earned	(5) Points Used
Post-merchant entry	-2.79*** (0.07)	-22.56*** (0.85)	-0.95*** (0.02)	-176.56*** (3.74)	-162.86*** (3.33)
Clothing retailer customers	-20.94*** (0.14)	-216.50*** (1.06)	-5.32*** (0.06)	-758.59*** (26.30)	-788.43*** (26.06)
Post-merchant entry * Clothing retailer customers	2.90*** (0.13)	19.79*** (1.19)	0.98*** (0.05)	294.08*** (31.77)	282.62*** (32.64)
Constant	15.19*** (0.11)	175.59*** (1.04)	3.58*** (0.04)	329.07*** (6.4680)	317.45*** (5.70)
Observations	475,081	475,081	475,081	279,879	279,879
R-squared	0.39	0.31	0.41	0.31	0.36
Day-of-week FE	YES	YES	YES	YES	YES

Store FE	YES	YES	YES	YES	YES
----------	-----	-----	-----	-----	-----

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

Table C-14 Shoes and accessories retailer entry, program-level analysis

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (in 000s)	(4) Points Earned	(5) Points Used
Post-merchant entry	-1,816.70*** (394.58)	-60.90*** (9.57)	-1,045.14*** (124.54)	-17,472.55 (24,041.10)	19,756.69 (13,974.40)
Shoe retailer customers	-8,472.44*** (380.18)	32.63** (14.15)	-3,318.34*** (113.54)	-151370.46*** (22,851.99)	-125641.70*** (12,193.97)
Post-merchant entry * Shoe retailer customers	1,887.07*** (421.18)	-62.59*** (17.49)	1,056.26*** (135.50)	19,617.77 (24,475.25)	-16,698.74 (14,413.53)
Constant	8,563.67*** (365.42)	405.29*** (6.54)	3,358.53*** (106.95)	152,960.85*** (22,614.76)	126,776.92*** (11,925.15)
Observations	362	362	362	362	362
R-squared	0.80	0.28	0.84	0.31	0.52
Day-of-week FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

Table C-15 Clothing and accessories retailer entry, store-level analysis

VARIABLES	(1) Total no. of transactions	(2) Average basket size	(3) Aggregate sales (in 000s)	(4) Points Earned	(5) Points Used
Post-merchant entry	-6.39*** (0.12)	-113.61*** (2.58)	-2.11*** (0.06)	-34.16*** (9.43)	43.15*** (6.85)
Shoe retailer customers	-14.86*** (0.14)	-340.07*** (2.30)	-5.78*** (0.06)	-638.70*** (19.44)	-598.55*** (14.44)
Post-merchant entry * Shoe retailer customers	6.92*** (0.15)	121.16*** (2.97)	2.38*** (0.08)	68.73*** (17.02)	11.93 (15.48)
Constant	16.21*** (0.19)	401.05*** (3.04)	6.12*** (0.07)	485.76*** (17.88)	402.43*** (10.90)
Observations	218,770	218,770	218,770	96,780	96,780
R-squared	0.40	0.27	0.42	0.11	0.18
Day-of-week FE	YES	YES	YES	YES	YES
Store FE	YES	YES	YES	YES	YES

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

C.4 Market Level Bayesian Structural Time Series

We replicate the BSTS modeling strategy described in the paper but apply to the full dataset excluding only the new merchant's transaction data. We interpret these findings as the national treatment effect on all stores and customers given the new merchant addition to the coalition loyalty program.

Our first measure is the average basket size of purchases at pre-existing stores (Figure C-1 below). We note that this series increases towards the end of the year shopping season, and decreases later in the observation period. In spite of the large basket spikes around the holiday and New Year's shopping period, the observed basket size shrinks well below the predictions of the counterfactual. We estimate that there is a posterior probability of 87% that the average basket size decreases by 60 of the local currency with a 95% Bayesian credible interval of [-189, 26]. This is equivalent to 12% decrease in the observed basket size with a 95% Bayesian credible interval of [-39%, 5.4%].

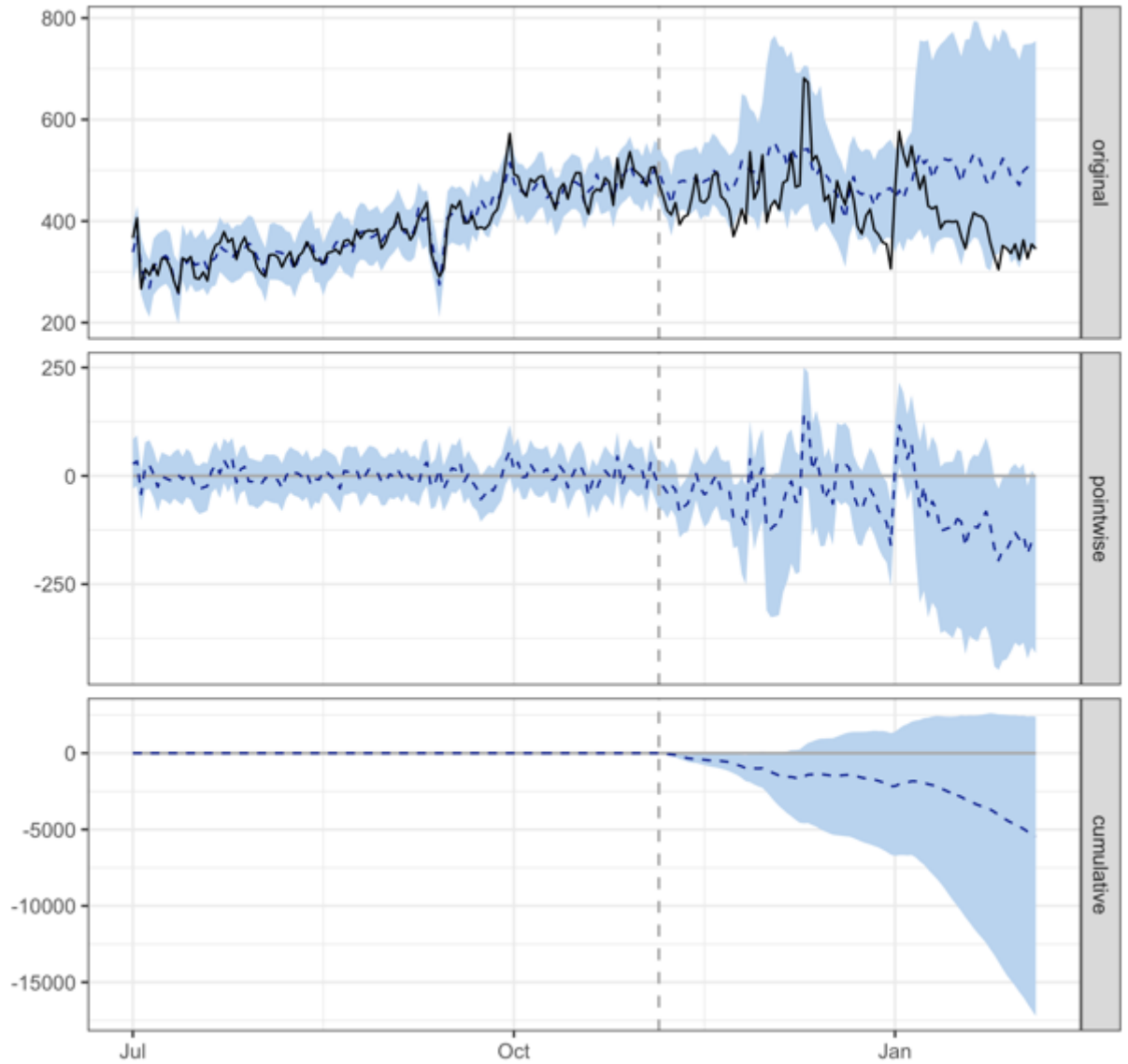


Figure C-1 Comparing observed data to counterfactual data for average basket size at pre-existing stores for the full data set

Our second measure is aggregate sales at pre-existing stores per day (Figure C-2 below). We note that this series increases steadily throughout the observation period, peaks in the holiday season and dips in January. The series moves similarly with the counterfactual indicating that there is no discernible change in aggregate sales for pre-existing stores in the three months following

their entry into the network. This is not too surprising as the small effect size from the matched sample regression model would be difficult to tease out in the aggregate data.

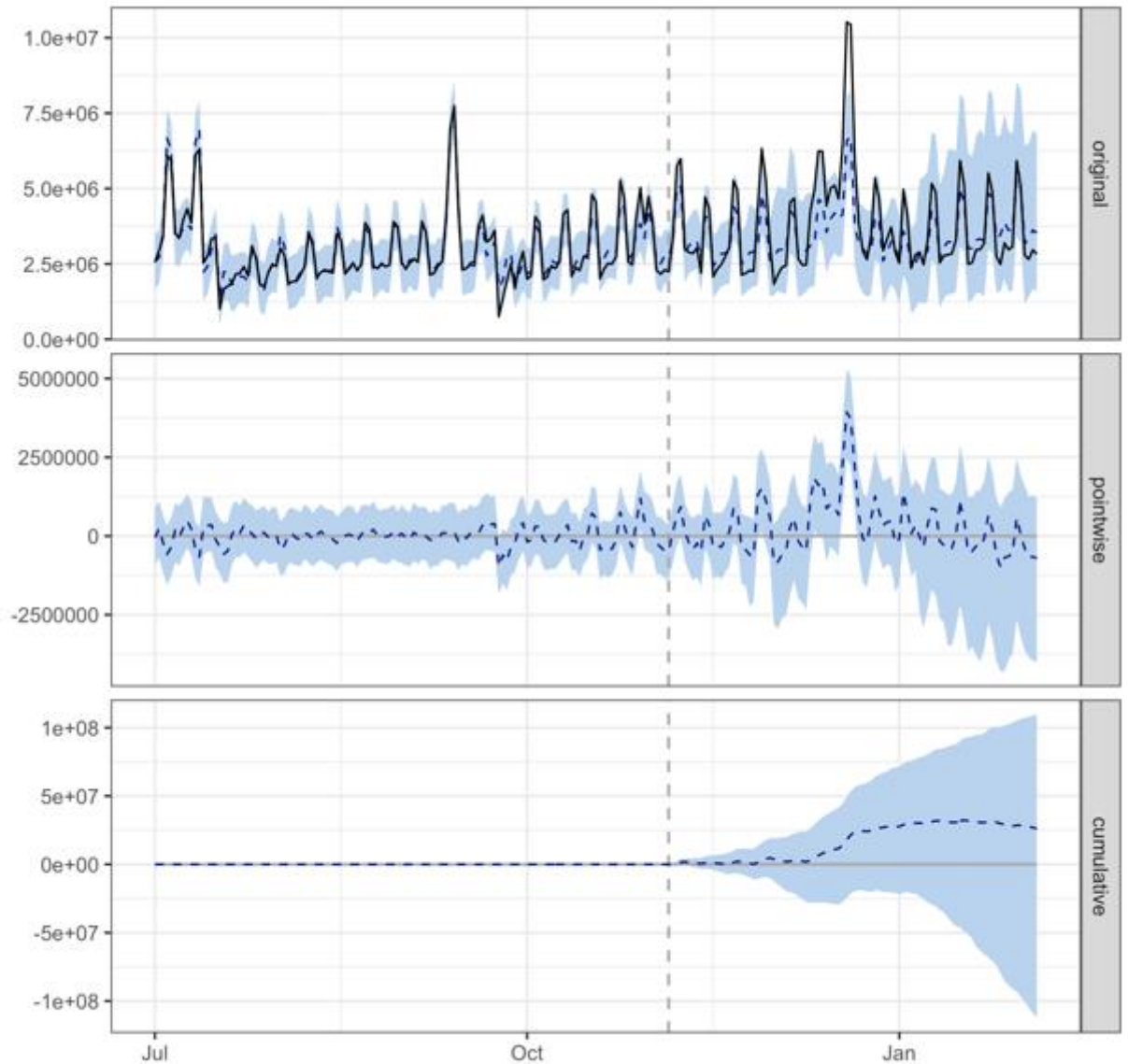


Figure C-2 Comparing observed data to counterfactual data for aggregate sales at pre-existing stores for the full data set

Our third measure is the total number of purchases at pre-existing stores per day (Figure C-3 below). We note that this series also increases steadily through the observation period, but that it increases faster than expected after the new merchant enters the loyalty program. This indicates

that consumers seem to be buying more frequently at pre-existing stores after the entry. We estimate that the number of daily purchases increased by 1,332 transactions on average. The posterior probability of a positive causal effect is 86%. The effect has a 95% Bayesian credible interval of [-1464, 3,880]. This is equivalent to 18% increase in the observed daily purchases with a 95% Bayesian credible interval of [-19%, 51%].

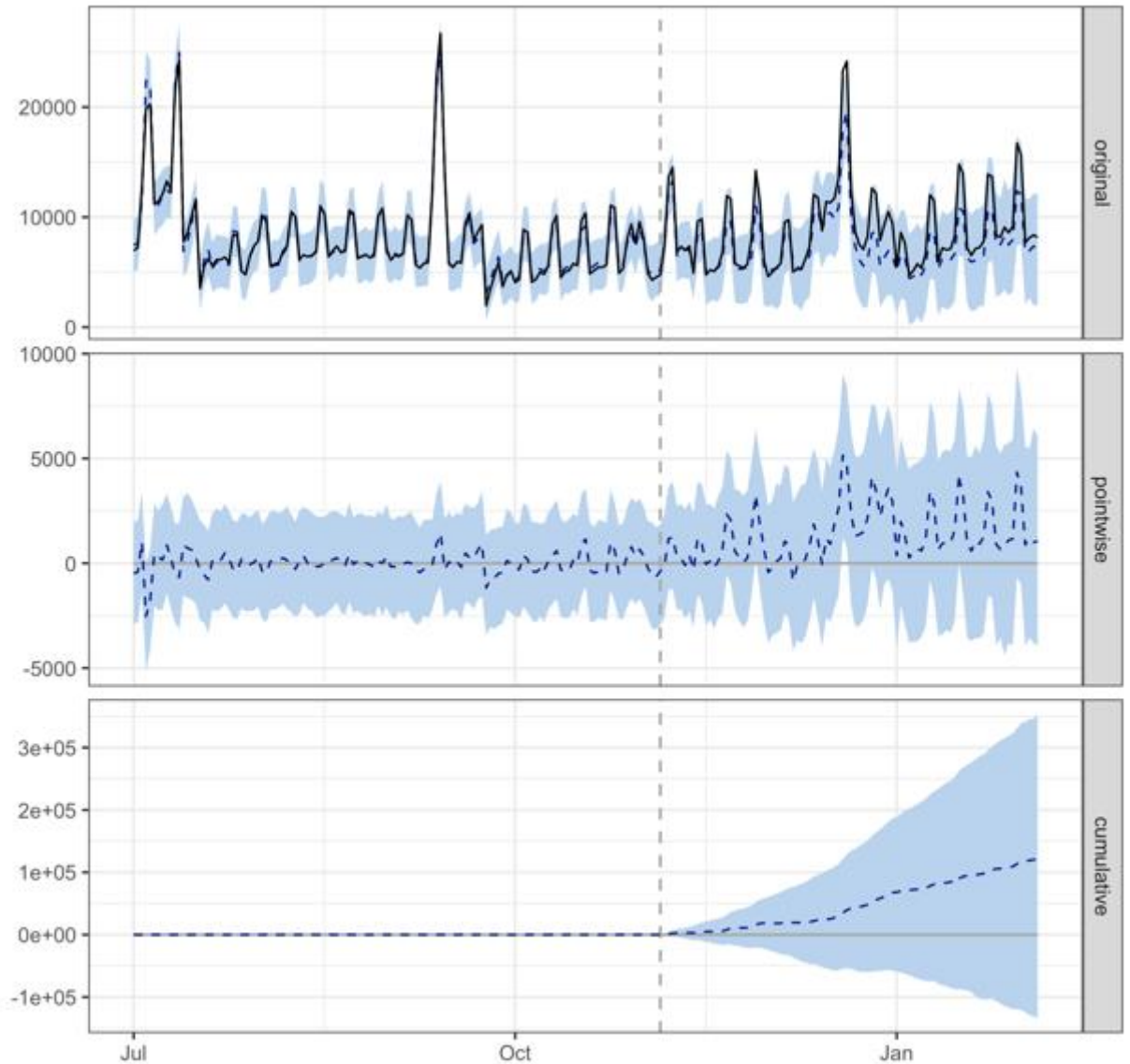


Figure C-3 Comparing observed data to counterfactual data for total no. of transactions at pre-existing stores for the full data set

Our fourth measure is the total number of loyalty points used at pre-existing stores per day (Figure C-4 below). We note that this series also increases steadily through the observation period, but that it increases faster than expected after the new merchant enters the loyalty program. This indicates that consumers seem to be using the loyalty points more frequently at pre-existing stores after the entry. We estimate that the number of daily points used increased by 92,697 points on

average. The posterior probability of a positive causal effect is 99.7%. The effect has a 95% Bayesian credible interval of [28,000, 160,000]. This is equivalent to a 67% increase in the observed daily purchases with a 95% Bayesian credible interval of [21%, 118%].

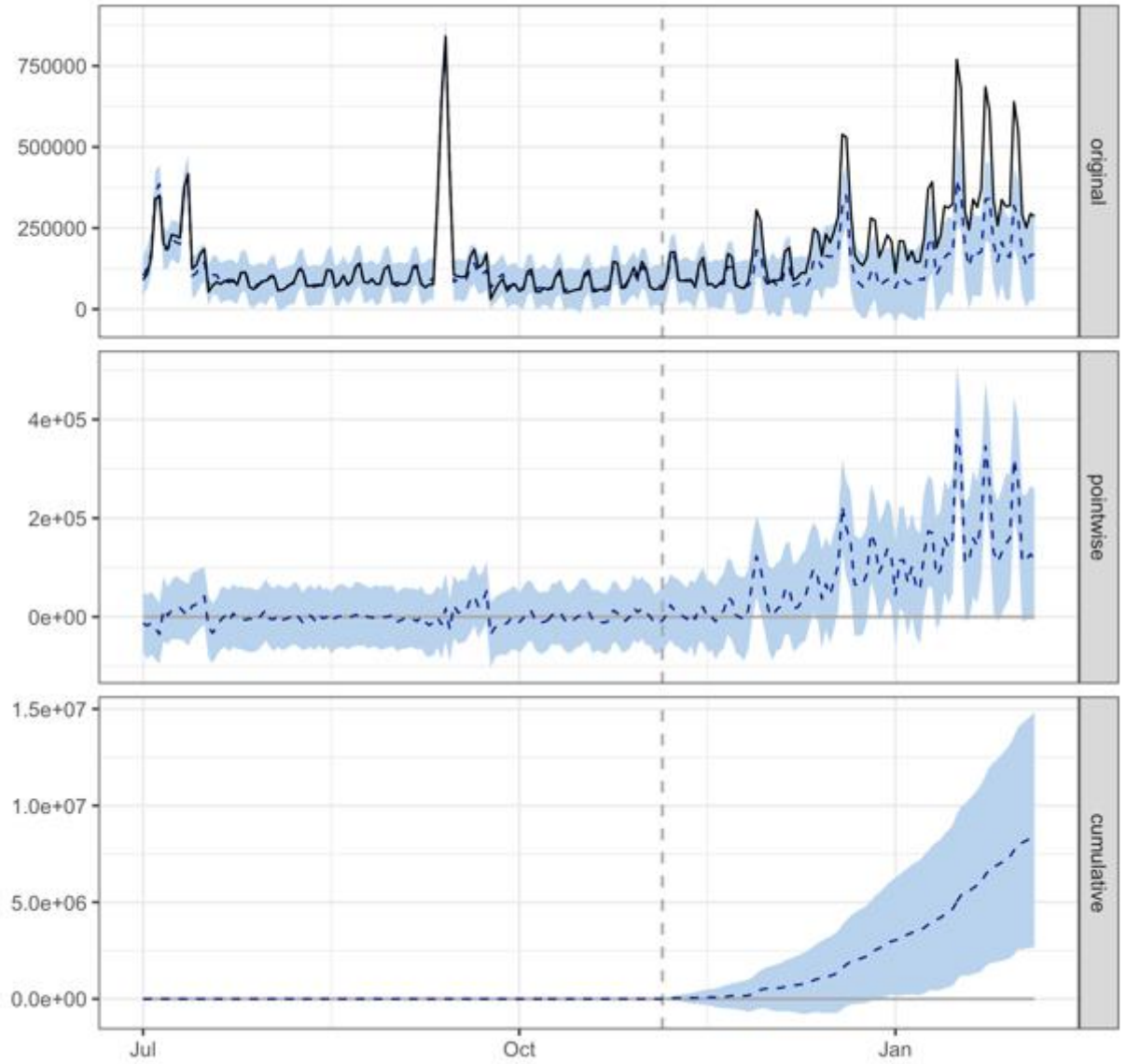


Figure C-4 Comparing observed data to counterfactual data for points used at pre-existing stores for the full data set

Our fifth measure is the total number of points earned at pre-existing stores per day (Figure C-4 below). We note that this series also increases steadily through the observation period, but not

more than expected after the new merchant enters the loyalty program. This indicates that consumers seem to be earning the loyalty points at the normal rate at pre-existing stores after the entry.

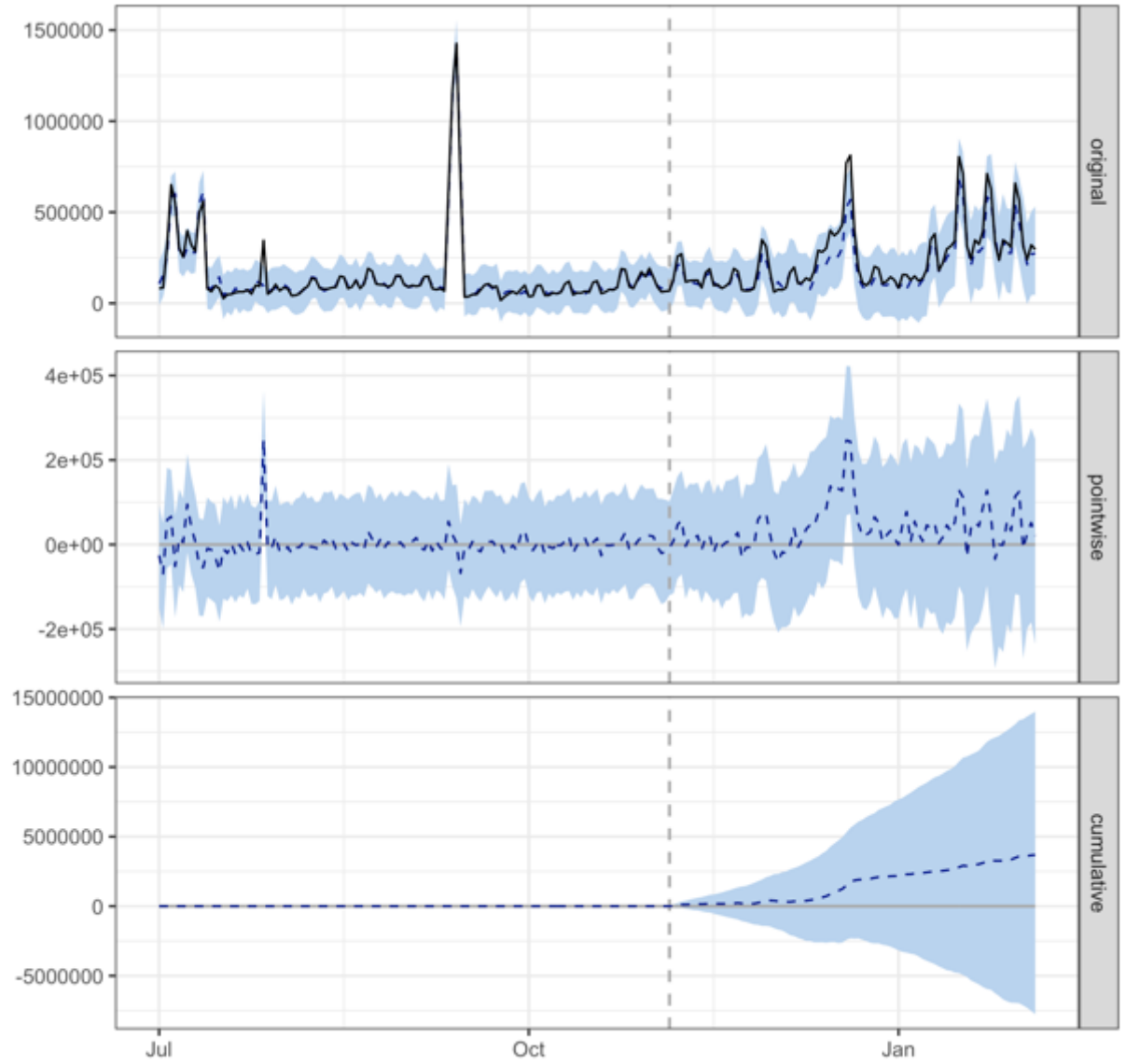


Figure C-5 Comparing observed data to counterfactual data for points earned at pre-existing stores for the full data set