

# On the Near Impossibility of Measuring the Returns to Advertising\*

Randall A. Lewis  
Google, Inc.  
ralewis@google.com

Justin M. Rao  
Microsoft Research  
justin.rao@microsoft.com

December 11, 2013

## Abstract

Classical theories assume the firm has access to reliable signals to measure the causal impact of choice variables on profit. For advertising expenditure we show, using twenty-five online field experiments with major U.S. retailers and brokerages (\$2.8 million expenditure), that this assumption typically does not hold. Evidence from the randomized trials is very weak because individual-level sales are incredibly volatile relative to the per capita cost of a campaign—a “small” impact on a noisy dependent variable can generate positive returns. A calibrated statistical argument shows that the required sample size for an experiment to generate informative confidence intervals is typically in excess of ten million person-weeks. This also implies that selection bias unaccounted for by observational methods only needs to explain a tiny fraction of sales variation to severely bias observational estimates. We discuss how weak informational feedback has shaped the current marketplace and the impact of technological advances moving forward.

**Keywords:** *advertising, field experiments, causal inference, electronic commerce, return on investment, information*

**JEL Codes:** *L10, M37, C93*

---

\*Previous versions circulated under the name “On the Near Impossibility of Measuring Advertising Effectiveness.” We especially thank David Reiley for his contributions to this work. Ned Augenblick, Arun Chandrasekhar, Sharad Goel, Garrett Johnson, Clara Lewis, R. Preston McAfee, Markus Möbius, Lars Lefgren, Michael Schwarz and Ken Wilbur gave us valuable feedback as well. We also thank attendees at Brigham Young University’s Economics Seminar, the Becker Friedman Institute Advances with Field Experiments Conference, Wharton OPIM, and other venues where we have presented this work. We also thank countless engineers, sales people, and product managers at Yahoo, Inc. Much of this work was done when the authors were at Yahoo! Research, Santa Clara, CA. The work represents our views and not those of our current or former employers.

# 1 Introduction

Each day a typical American sees 25–45 minutes of television commercials, many billboards and numerous online ads (Kantar Media, 2008). Industry reports place annual U.S. advertising revenue in the range of \$173 billion, or about \$500 per American per year. This means that in order to break even, the universe of advertisers needs to net roughly \$1.50 in profits per person per day, corresponding to about \$4–6 in incremental sales per person per day, or about \$3,500–5,500 per household per year. These back-of-the-envelope calculations demonstrate that the market valuation of advertising implies a large causal impact on household purchases in order for advertising to be profitable on average, yet, perhaps surprisingly, it is an open question as to whether this is in fact the case.

Consistent with this aggregate uncertainty, papers in the advertising effectiveness literature often use “Do ads have any effect?” as the key hypothesis to test. This is perhaps epitomized by the *first sentence* of an influential paper by Abraham and Lodish (1990), “Until recently, believing in the effectiveness of advertising and promotion was largely a matter of faith.” An opening that might seem a bit peculiar, given that before it was written American firms had spent approximately \$4.6 trillion promoting their products and services.

In this paper we address the underlying puzzle: if so much money is being spent on advertising, how could it be possible that firms have such imprecise beliefs on the returns? We present a data-driven argument that shows precisely why this is this case and, furthermore, why it will remain the case for the foreseeable future. It turns out that a key assumption of the classical theory of the firm, namely access to reliable signals mapping choice variables to profit, tends to fail in this domain. This assertion is based on our analysis of 25 large-scale digital advertising *field experiments* from well-known retailers and financial service firms partnering with a large web publisher. In total, the experiments accounted for \$2.8 million in expenditure. We find that even when ad delivery and consumer purchases can be measured at the individual level, linked across purchasing domains, and randomized to ensure exogenous exposure, forming reliable estimates on the returns to advertising is exceedingly difficult. As an advertiser, the data are stacked against you.

The intuition for the inference difficulty can be gleaned from the following observation: the effect of a single campaign should be “small” in equilibrium. Most

ads are relatively cheap (typically  $< \$0.02$  per delivery) so only a small fraction of people need to be “converted” for a campaign to be profitable. Using detailed sales data from our partner firms, we show that matters are further complicated by the fact that the standard deviation of sales, on the individual level, is typically ten times the mean over the duration of a typical campaign and evaluation window. (While this relationship may not hold true for smaller firms or new products, it was remarkably consistent across the relatively diverse set of large advertisers in our study.) As a consequence, the advertiser has to estimate a relatively subtle effect in an incredibly noisy economic environment.

To provide a framework for our empirical analysis we develop a simple model of the firm’s advertising problem. The key model parameter is return on investment (ROI)—the profits generated through the advertising as a percentage of the costs. In online advertising, intermediate metrics such as clicks have become popular in measuring advertising effectiveness.<sup>1</sup> Our focus, however, is on what the firm presumably cares about in the end, namely profits, and our extensive data sharing agreements allow us to sidestep intermediate metrics. We show that even a fully randomized experiment, massive trials (typically in the single-digit millions of person-weeks at typical advertising intensity) are required to reliably distinguish disparate hypotheses such as “the campaign had no effect” ( $-100\%$  ROI) from “the campaign was profitable for the firm” ( $\text{ROI} > 0\%$ ). Answering questions such as “was the ROI 15% or -5%,” a large difference for your average investment decision, or “was the *annualized* ROI at least 5%,” a reasonable question to calibrate against the cost of capital, typically requires at least hundreds of millions of independent person-weeks—nearly impossible for a campaign of any realistic size. And while ROI measures accounting profits and losses, determining the profit-maximizing level of ROI requires one to estimate the underlying profit function. We briefly discuss the (rather incredible) difficulties of this enterprise.

Rather than endorsing “big data” observational methods, the shortcomings of experiments instead serve to highlight their lurking biases. Ads are, by design, not delivered randomly because marketers target across time, people, and contexts. So while the true causal effect should be relatively small, selection effects are expected to be quite large. Consider a simple example: if an ad costs 0.5 cents per delivery

---

<sup>1</sup>For a discussion of complications that can arise from using these metrics, see Lewis, Rao, and Reiley (2013).

(typical of “premium” online display ads), each viewer sees one ad, and the marginal profit per “conversion” is \$30, then only 1 in 6,000 people need to be “converted” for the ad to break even. Suppose a targeted individual has a 10% higher baseline purchase probability (a very modest degree of targeting), then the selection effect is expected to be *600 times larger* than the causal effect of the ad.

We present our results in considerable detail in Section 3, but the core intuition can be gleaned from a high-level look at a representative (median) campaign. Imagine running a regression of sales per individual (in dollars) on an indicator variable of whether or not the person saw advertising. In an experiment, the indicator variable is totally exogenous, while in an observational method one attempts to control for selection bias. Given the approximately \$100,000 cost and reach in the low millions of people of our representative campaign, to net a +25% ROI, the campaign must causally raise sales by \$0.35 per person. Incorporating the volatility of sales reveals that this regression amounts to detecting 35 cent impact on a variable with a mean of \$7 and a standard deviation of \$75. This implies that the  $R^2$  for a *highly profitable* campaign is on the order of 0.0000054.<sup>2</sup> To successfully employ an observational method, we must be sure we have not omitted any control variables or misspecified the functional form to a degree that would generate an  $R^2$  on the order of 0.000002 or more, otherwise estimates will be *severely* biased. This seems to be an impossible feat to accomplish in any circumstance but especially in one where selection effects are expected to be orders of magnitude larger than the true causal effect.

Since we are making the (admittedly) strong claim that most advertisers do not, and indeed many *cannot*, know the effectiveness of their advertising spend, it is paramount to stress test the generalizability of the empirical results this claim rests on. We first show that the firms we study are fairly representative in terms of sales volatility, margins, and size of advertisers that constitute the majority of ad spending. A caveat is that our results do not necessarily apply to small firms or new products. Our tongue-in-cheek “Super Bowl ‘Impossibility’ Theorem” shows that even a massive, idealized experiment would be relatively uninformative for many advertisers. We theoretically demonstrate that our results were not driven by the campaign windows we chose or the level of targeting of the campaigns under study. Finally, we discuss how recent empirical work from a major advertiser’s research lab

---

<sup>2</sup> $R^2 = \frac{1}{4} \cdot \left(\frac{\$0.35}{\$75}\right)^2 = 0.0000054.$

helps buttress our central claim (Blake et al., 2013).

Our central claim in turn has deep implications for the advertising and publishing market. Scarce information means there is little “selective pressure” on advertising levels across firms. Consistent with this reasoning, we use additional data to show that otherwise similar firms in the same industries having vastly different levels of advertising spending. As experimentation becomes more common, the informational landscape will increasingly shape the web publishing market by setting massive firms off to an advantage because their scale is necessary for experiments to provide reliable feedback on the returns to advertising.

The paper proceeds as follows: Section 2 gives a simple model of the advertiser’s problem and calibrates it using empirical values from our study, Section 3 presents the main empirical results, and in Section 4 we discuss these results in the context of the broader market. We briefly conclude in Section 5.

## 2 A Simple Model of the Advertiser’s Problem

In this section we formalize the problem of campaign evaluation. Our model is exceedingly simple, designed to capture the core elements of measuring advertising returns.

### 2.1 Model

Full-blown optimization of advertising would include, among other things, selecting the consumers to advertise to, measuring the advertising technology across various media, and determining how those technologies interact. Our focus here is on measuring the returns to an advertising campaign—the crucial first step in optimization. We define a campaign as a set of advertisements delivered to a set of consumers through a single channel over a specified (and typically short) period of time using one “creative” (all messaging content such as pictures, text, and audio). Ex-post evaluation asks the question, “Given a certain expenditure and delivery of ads, what is the rate of return on this investment (ROI)?” Note that for now we take the target population as given.

A campaign is defined by  $c$ , the cost per user. For a given publishing channel,  $c$  determines how many “impressions” each user sees. We assume the sales impact is

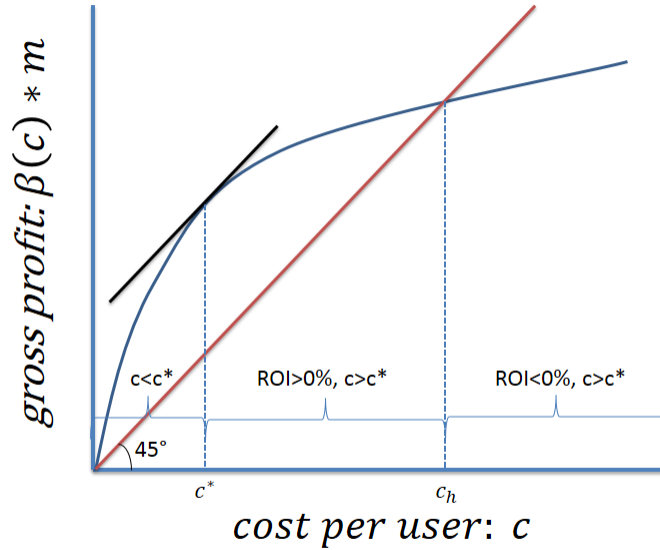


Figure 1: Graphical depiction of the advertiser’s problem.

defined by a continuous concave function of per-user expenditure  $\beta(c)$ .<sup>3</sup> We can easily incorporate consumer heterogeneity with a mean-zero multiplicative parameter on this function and then integrate this parameter out to focus on the representative consumer. Let  $m$  be the gross margin of the firm so that  $\beta(c) * m$  gives gross profit per person. Net profit subtracts cost  $\beta(c) * m - c$ , and ROI measures net profit as a percentage of cost  $\frac{\beta(c)*m-c}{c}$ . In our simple model the only choice variable is  $c$ , or “how much I advertise to each consumer.”

Figure 1 graphically depicts the model:  $c^*$  gives optimal spend and  $c_h$  gives the spend level where ROI is exactly 0%. At any point past  $c_h$  the firm has negative returns, whereas any point to the left of  $c^*$  the firm has positive returns but is under-advertising. For points in  $(c^*, c_h)$ , the firm is over-advertising because marginal return is negative but average return, or ROI, is still positive.

The model formalizes the estimation of the average per-person impact of a given

<sup>3</sup>For supportive evidence of concavity see (Lewis, 2010). This assumption could be weakened to “concave in the region of current spending,” which essentially just says that the returns to advertising are not infinite and the firm is not in a convex region.

campaign on consumer behavior. In reality, multiple creatives are used, the actual quantity of ads delivered per person is stochastic (because exposure depends on user activity), and  $\beta$  would include arguments such as the spending on other advertising channels. Our evaluation framework is motivated by the fact that the “campaign” is an important operational unit in marketing. A Google Scholar search of the exact phrase “advertising campaign” returned 48,691 unique research documents. This is echoed by our personal experience as well.

## 2.2 Measuring the Returns to a Campaign

We start out with a high-level view of the inference challenges facing an advertiser by calibrating the model using median figures from our experiments. On the cost side, online display ad campaigns that deliver a few ads per day per person cost about 1–2 cents per person per day and typically run for about two weeks, cumulating in a cost between 15 and 40 cents per person (roughly the cost of one 30-second TV ad per person per day). Given the total volume of advertising a typical consumer sees across all media, even an intense campaign only captures about 2% of a targeted person’s advertising “attention.”<sup>4</sup>

Sales volatility has three components: the average magnitude (mean sales), heterogeneity (variance of per-person means), and rarity of purchases (stochasticity in purchasing). For the large retailers and financial service firms in our study, the mean weekly sales per-person varied considerably across firms, as does the standard deviation in sales. However, we find that the ratio of the standard deviation to the mean (the coefficient of variation of the mean) is typically around ten for the retail firms—customers buy goods relatively infrequently, but when they do, the purchases tend to be quite large relative to the mean.<sup>5</sup> Sales volatility tends to be higher for financial service firms, because people either sign-up and become lucrative long-term customers or do not use the service at all.

---

<sup>4</sup>Ads are typically sold by delivered impressions, but this does not necessarily mean a person noticed them. It is possible for a campaign to get 100% of a consumer’s attention (he or she pays attention to that ad and ignores all others) or 0% (it is totally ignored) or any value in between.

<sup>5</sup>An extreme example of this feature is automobiles (which we discuss later) where the sales impact is either a number ranging in the tens of thousands of dollars, or more likely, given the infrequency of car purchases, it is \$0. Homogeneous food stuffs have more stable expenditure, but their very homogeneity likely reduces own-firm returns to and equilibrium levels of advertising within industry as a result of positive advertising spillovers to competitor firms (Kaiser, 2005).

In the econometric specification let  $y_i$  be sales for individual  $i$ . Since we are assuming, for simplicity, that each affected individual saw the same value of advertising for a given campaign, we will use an indicator variable  $x_i$  to quantify ad exposure.  $\hat{\beta}(c)$  gives our estimate of the sales impact for a campaign of cost-per-user  $c$ . Standard econometric techniques estimate this value using the difference between the exposed (E) and unexposed (U) groups. In an experiment, exposure is exogenous. In an observational study, one would also condition on covariates  $W$  and a specific functional form, which could include individual fixed effects, and the following notation would use  $y|W$ . All the following results go through with the usual “conditional upon” caveat.

For the case of a fully randomized experiment, our estimation equation is simply:

$$y_i = \beta x_i + \epsilon_i \tag{1}$$

We suppress  $c$  in the notation because a given campaign has a fixed size per user. The average sales impact estimate,  $\hat{\beta}$ , can be converted to ROI by multiplying by the gross margin to get the gross profit impact, subtracting per-person cost, and then dividing by cost to get the percentage return. In our empirical analysis we condition on all available covariates, such as lagged sales, to soak up residual variation; the arguments in this section are not affected by ignoring this strategy for now.

Below we use standard notation to represent the sample means and variances of the sales of the exposed and unexposed groups, the difference in means between those groups, and the estimated standard error of that difference in means. Without loss of generality we assume that the exposed and unexposed samples are the same size ( $N_E = N_U = N$ ) and have equal variances ( $\sigma_E = \sigma_U = \sigma$ ), which is the best-case scenario from a design perspective.

$$\bar{y}_E \equiv \frac{1}{N_E} \sum_{i \in E} y_i, \bar{y}_U \equiv \frac{1}{N_U} \sum_{i \in U} y_i \tag{2}$$

$$\hat{\sigma}_E^2 \equiv \frac{1}{N_E - 1} \sum_{i \in E} (y_i - \bar{y}_E)^2, \hat{\sigma}_U^2 \equiv \frac{1}{N_U - 1} \sum_{i \in U} (y_i - \bar{y}_U)^2 \tag{3}$$

$$\Delta \bar{y} \equiv \bar{y}_E - \bar{y}_U \tag{4}$$



$$\hat{\sigma}_{\Delta\bar{y}} \equiv \sqrt{\frac{\hat{\sigma}_E^2}{N_E} + \frac{\hat{\sigma}_U^2}{N_U}} = \sqrt{\frac{2}{N}} \cdot \hat{\sigma} \quad (5)$$

We focus on two familiar econometric statistics. The first is the  $R^2$  of the regression of  $y$  on  $x$ , which gives the fraction of the variance in sales attributed to the campaign (or, in the model with covariates, the partial  $R^2$  after first conditioning on covariates in a first stage regression—for a nice explanation of how this works, see Lovell, 2008):

$$R^2 = \frac{\sum_{i \in U} (\bar{y}_U - \bar{y})^2 + \sum_{i \in E} (\bar{y}_E - \bar{y})^2}{\sum_i (y_i - \bar{y})^2} = \frac{2N \left(\frac{1}{2}\Delta\bar{y}\right)^2}{2N\hat{\sigma}^2} = \frac{1}{4} \left(\frac{\Delta\bar{y}}{\hat{\sigma}}\right)^2. \quad (6)$$

Second is the  $t$ -statistic for testing the hypothesis that the advertising had no impact ( $\beta = 0$ ):

$$t_{\Delta\bar{y}} = \frac{\Delta\bar{y}}{\hat{\sigma}_{\Delta\bar{y}}} = \sqrt{\frac{N}{2}} \left(\frac{\Delta\bar{y}}{\hat{\sigma}}\right). \quad (7)$$

In both cases, we have related a standard regression statistic to the ratio between the average impact on sales ( $\Delta\bar{y}$ ) and the standard deviation of sales ( $\sigma$ )—we will call this the “impact-to-standard-deviation ratio.” It is also known as *Cohen’s d*.

We calibrate the test statistics using median values from 19 experiments run with large U.S. retailers in partnership with Yahoo! (the remaining 6 experiments were for account sign-ups for financial firms, making it harder to determine sales in dollars). For ease of exposition, we will discuss the hypothetical case as if it were a single, actual experiment. This representative campaign costs \$0.14 per customer, which amounts to delivering 20–100 display ads at a price of \$1–\$5 CPM,<sup>6</sup> and the gross margin is assumed to be 50%.<sup>7</sup> Mean sales per-person for the period under study is \$7 and the standard deviation is \$75.

We will suppose the ROI goal was 25%, which corresponds to a \$0.35 sales impact per person, yielding gross profits of \$0.175 per person as compared to costs of \$0.14. A \$0.35 per-person impact on sales corresponds to a 5% increase in sales during the two weeks of the campaign (note that in terms of percentage lift, the required impact of the campaign appears quite large). The estimation challenge

---

<sup>6</sup>CPM is the standard for impression-based pricing for online display advertising. It stands for “cost per mille” or “cost per thousand.”

<sup>7</sup>We base this assumption on our conversations with retailers, our knowledge of the industry and SEC filings.

facing the advertiser is to detect this \$0.35 difference in sales between the treatment and control groups amid the noise of a \$75 standard deviation in sales. The impact-to-standard-deviation ratio is only 0.0047. From our derivation above, this implies an  $R^2$  of:

$$R^2 = \frac{1}{4} \cdot \left( \frac{\$0.35}{\$75} \right)^2 = 0.0000054. \quad (8)$$

Even a very successful campaign has a  $R^2$  of only  $0.0000054$ , meaning we need a very large  $N$  to reliably distinguish it from 0, let alone give a precise confidence interval. Suppose we had 2 million unique users evenly split between test and control in a fully randomized experiment. With a true ROI of 25% and an impact-to-standard-deviation ratio of 0.0047, the expected  $t$ -statistic with a null hypothesis of 100% ROI (zero causal impact) is 3.30, using the above formula. This corresponds to a test with power of about 95% at the 10% (5% one-sided) significance level because the approximately normally distributed  $t$ -statistic should be less than the critical value of 1.65 about 5% of the time (corresponding to the cases where we cannot reject the null). With 200,000 unique users, the expected  $t$ -statistic is 1.04, indicating an experiment of this size is hopelessly underpowered: under the alternative hypothesis of a healthy 25% ROI, we fail to reject the null that the ad had no causal impact 74% of the time.<sup>8</sup>

The minuscule  $R^2$  for the treatment variable in our representative randomized trial has serious implications for observational studies, such as regression with controls, difference-in-differences, and propensity score matching. An omitted variable, misspecified functional form, or slight amount of correlation between browsing behavior and sales behavior generating  $R^2$  on the order of 0.0001 is a *full order of magnitude* larger than the true treatment effect. Meaning a very small amount of endogeneity would *severely bias* estimates of advertising effectiveness. Compare this to a classic economic example of wage/schooling regressions, in which the endogeneity is typically 1/8 the treatment effect (Card, 1999). If the partial  $R^2$  of the treatment variable is very small, clean identification becomes paramount. As we showed in the introduction, a minimal level of targeting that results in the exposed group having a 10% higher baseline purchase rate can lead to an exposed-unexposed

---

<sup>8</sup>Note that when a low powered test does, in fact, correctly reject the null, the point estimates conditional on rejecting will be significantly larger than the alternatively hypothesized ROI. That is, when one rejects the null, the residual on the estimated effect is positive. This overestimation was recently dubbed the “exaggeration factor” by Gelman and Carlin (2013).

difference of about 600 times the true treatment effect. Unless this difference is controlled for with near *perfect* precision, observational models will have a large bias.

In demonstrating these lurking biases, are we arguing against a straw man? Not so, according to a recent article in the *Harvard Business Review*. The following quotation is from the president of comScore, a large data-provider for web publishers and advertisers:

Measuring the online sales impact of an online ad or a paid-search campaign—in which a company pays to have its link appear at the top of a page of search results—is straightforward: We determine who has viewed the ad, then compare online purchases made by those who have and those who have not seen it.

M. Abraham, 2008, *Harvard Business Review*.

The author used this methodology to report a 300% improvement in outcomes for the exposed group, which seems surprisingly high as it implies that advertising prices should be at least an order of magnitude higher than current levels.

## 3 Analysis of the 25 Field Experiments

### 3.1 Summary Statistics and Overview

Table 1 gives an overview of 25 display advertising experiments/campaigns. We highlight the most important figures and present summary statistics. Due to confidentiality agreements, we cannot reveal the identity of the advertisers. We can say they are large firms that are most likely familiar to American readers. We employ a naming convention using the vertical sector of the advertiser in lieu of the actual firm names. The firms in Panel 1 are retailers, such as large department stores; in Panel 2 they are financial service firms.<sup>9</sup>

Columns 1–3 of Table 1 give basic descriptors of the experiment. Sales is the key dependent measure for the firms in Panel 1 and Column 4 gives the unit of observation (“3” indicates daily observation, “4” indicates weekly). In Panel 2, the

---

<sup>9</sup>Some of the experiments are taken from past work out of Yahoo! Labs, such Lewis and Reiley (2013).

Table 1: Overview of the 25 Advertising Field Experiments

Retailers: In-Store + Online Sales*														
Estimation Strategies Employed**														
Adv	Year	#	Y	X	Y&X	W	Days	Cost	Campaign Level Summary				Per Customer	
									Test	Control	Exposed	Control	Avg. Sales (Control)	$\sigma$ sales
R 1	2007	1	1,4	1	-	1,2,3	14	\$128,750	1,257,756	300,000	814,052	-	\$9.49	\$94.28
R 1	2007	2	1,4	1	-	1,2,3	10	\$40,234	1,257,756	300,000	686,878	-	\$10.50	\$111.15
R 1	2007	3	1,4	1	-	1,2,3	10	\$68,398	1,257,756	300,000	801,174	-	\$4.86	\$69.98
R 1	2008	1-6	1,4	1,2,3	-	1,2,3	105	\$260,000	957,706	300,000	764,235	238,904	\$125.74	\$490.28
R 1	2010	1	1,4	1,2	-	1,2,3,4	7	\$81,433	2,535,491	300,000	1,159,100	-	\$11.47	\$111.37
R 1	2010	2-3	1,3,4	1,2,3,4	1	1,2,3	14	\$150,000	2,175,855	1,087,924	1,212,042	604,789	\$17.62	\$132.15
R 2	2009	1a	1,5	1	-	-	35	\$191,750	3,145,790	3,146,420	2,229,959	-	\$30.77	\$147.37
R 2	2009	1b	1,5	1	-	-	35	\$191,750	3,146,347	3,146,420	2,258,672	-	\$30.77	\$147.37
R 2	2009	1c	1,5	1	-	-	35	\$191,750	3,145,996	3,146,420	2,245,196	-	\$30.77	\$147.37
R 3	2010	1	1,3,4	1,2,3	1	1,3	3	\$9,964	281,802	161,163	281,802	161,163	\$1.27	\$18.46
R 3	2010	2	1,3,4	1,2	1	1,3	4	\$16,549	483,015	277,751	424,380	-	\$1.08	\$14.73
R 3	2010	3	1,3,4	1,2,3	1	1,3	2	\$25,571	292,459	169,024	292,459	169,024	\$1.89	\$18.89
R 3	2010	4	1,3,4	1,2,3	1	1,3	3	\$18,234	311,566	179,709	311,566	179,709	\$1.29	\$16.27
R 3	2010	5	1,3,4	1,2	1	1,3	3	\$18,042	259,903	452,983	259,903	-	\$1.75	\$18.60
R 3	2010	6	1,3,4	1,2,3	1	1,3	4	\$27,342	355,474	204,034	355,474	204,034	\$2.64	\$21.60
R 3	2010	7	1,3,4	1,2,3	1	1,3	2	\$33,840	314,318	182,223	314,318	182,223	\$0.59	\$9.77
R 4	2010	1	1,3,4	1,2	1	1	18	\$90,000	1,075,828	1,075,827	693,459	-	\$0.56	\$12.65
R 5	2010	1	1,5	1,2	-	1,3	41	\$180,000	2,321,606	244,432	1,583,991	-	\$54.77	\$170.41
R 5	2011	1	1,3,4	1,2	1	1,3	32	\$180,000	600,058	3,555,971	457,968	-	\$8.48	\$70.20

Financial Services: New Accounts Online Only\*\*\*

Estimation Strategies Employed**														
Adv	Year	#	Y	X	Y&X	W	Days	Cost	Campaign Level Summary				Per Customer	
									Test	Control	Exposed	Control	Pr New (Test)	SE New Acct
F 1	2008	1a	2,5	1,2,4	-	3	42	\$50,000	12% of Y!	52% of Y!	794,332	867	0.0011	0.0330
F 1	2008	1b	2,5	1,2,4	-	3	42	\$50,000	12% of Y!	52% of Y!	748,730	762	0.0010	0.0319
F 1	2008	1c	2,5	1,2,4	-	3	42	\$75,000	12% of Y!	52% of Y!	1,080,250	1,254	0.0012	0.0341
F 1	2008	1d	2,5	1,2,4	-	3	42	\$75,000	12% of Y!	52% of Y!	1,101,638	1,304	0.0012	0.0344
F 2	2009	1	2,3	1,2,3,4	1,2	3	42	\$612,693	90% of Y!	10% of Y!	17943572	10,263	0.0006	0.0239
F 2	2011	1	2,5	1,2	-	4	36	\$85,942	8,125,910	8,125,909	793,042	1090	0.0014	0.0331

\* These retailers do a supermajority of sales via their brick & mortar stores.

\*\* Estimation strategies employed to obtain the standard errors of the ad impact between the test and control groups follow:

“Y” 1:Sales, 2:Sign-ups, 3:Daily, 4:Weekly, 5:Total Campaign Window;

“X” 1:Randomized Control, 2:Active on Y! Network or site where ads were shown, 3:Placebo Campaign for Control Group, 4:Multiple Treatments;

“Y&X” 1: Sales filtered post first exposure or first page view, 2:Outcome filtered based on post-exposure time window);

“W” 1:Lagged sales, 2:Demographics, 3:Online behaviors.

\*\*\* These financial services advertisers do a supermajority of their business online.

dependent measure is new account sign-ups. Column 7 gives the control variables, such as lagged sales, we have to reduce noise in the experimental estimates. The experiments ranged from 2 to 135 days (Column 8), with a median of 14 days, which is typical of display campaigns. Column 9 shows the campaign cost varied from relatively small (\$9,964) to quite large (\$612,693). The mean was \$114,083; the median was \$75,000. The median campaign reached over one million individuals, and all campaigns had hundreds of thousands of individuals in both test and control cells (Columns 9–11). Overall, the campaigns represent over \$2.8 million in expenditure.

The second-to-last column shows that the average sales per customer varied widely across the firms. This is driven by both the popularity of the retailer and the targeting level of the campaign (a more targeted campaign typically has higher baseline sales). Median per person sales is \$8.48 for the test period. The final column gives the standard deviation of sales on an individual level. The median campaign had a standard deviation 9.83 times the mean. We plot the distribution of standard-deviation-to-mean ratio against campaign duration in Figure 2. This ratio exceeds seven for all but two of the experiments. Longer campaigns tend to have a lower ratio, which is due to sufficient independence in sales across weeks.<sup>10</sup> While a longer campaign (of the same size) generates more points of observation, these additional data will only make inference easier if the spending per person per week is not diluted (see section 4.1.4).

### 3.2 Estimating ROI

In Table 2, we take a detailed look at estimating ROI. Column 3 gives the standard error associated with the estimate of  $\beta$ , the test-control sales difference as defined by the model (in dollars for Panel 1, in account sign-ups for Panel 2). We condition on the control variables outlined in Column 7 of Table 1 in order to obtain as precise an estimate as possible. In Column 4, we give the implied radius (+/- window) of the 95% confidence interval for the sales impact, in percentage terms—the median radius is 5.5%. Column 5 gives the per-person advertising spend, which can be compared to the standard error of the treatment effect given in Column 3 to capture how the

---

<sup>10</sup>If sales are, in fact, independent across weeks, we would expect the coefficient of variation to follow  $\frac{\sqrt{T} \cdot \sigma_{weekly}}{T \cdot \mu}$ . However, over long horizons (i.e., quarters or years), individual-level sales are correlated, which also makes past sales a useful control variable when evaluating longer campaigns.

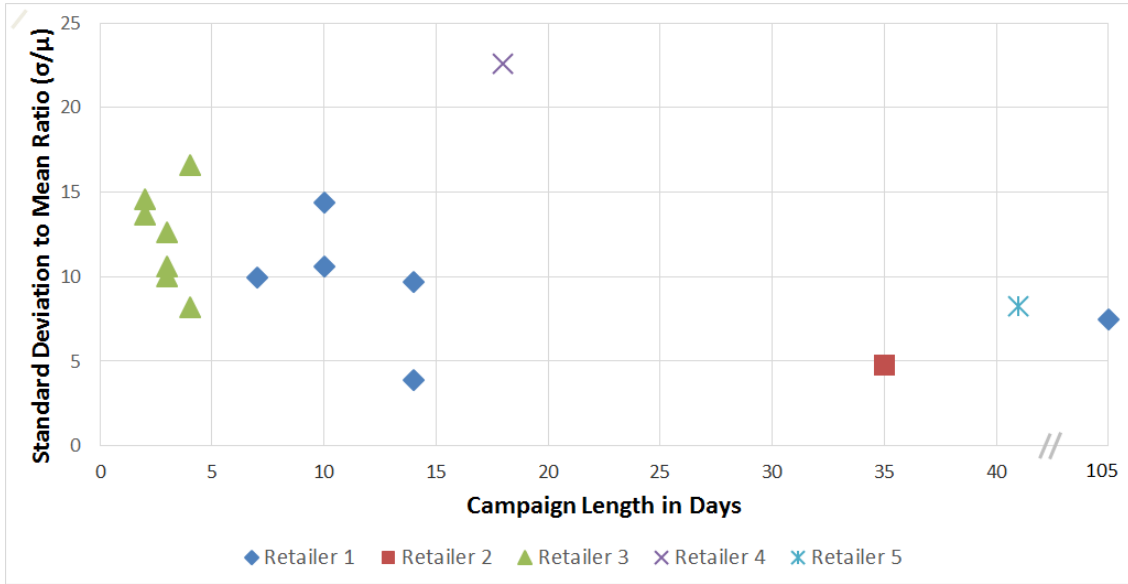


Figure 2: Relationship between sales volatility, as given by the coefficient of variation,  $\frac{\sigma}{\mu}$ , and campaign length in days.

magnitude of statistical uncertainty in sales relates to expenditure. In Column 7 we translate the sales impact standard errors to ROI using our estimates of gross margins (Column 6, based on SEC filings). For the financial firms we convert a customer acquisition into an approximate dollar value using figures they provided to us. The median standard error for ROI is 26.1%—it follows that the median confidence interval is about *100 percentage points wide*. The mean standard error is higher still at 61.8%, implying a confidence interval that is too wide to be of much practical use.

In Figure 3 we plot the standard error of the ROI estimate against the per capita campaign cost. Each line represents a different advertiser. Two important features are immediately apparent. First, there is significant heterogeneity across firms. Retailer 1 and the financial firms had the highest statistical uncertainty in the ROI estimate. Financial firms operate in an all-or-nothing world—someone either signs up for an account and likely becomes a lucrative long-term customer or does not, generating zero revenue. Retailer 1 simply had a higher standard deviation of sales. Second, estimation tends to get more precise as the per-person spend increases. The curves are downward sloping with the exception of a single point. This is exactly

Table 2: Statistical Precision and Power Calculations for the 25 Advertising Field Experiments

In-Store + Online Sales		Key Statistical Properties of Campaign										HARD			HARDER			HARDEST			CRAZY								
		Adv	#	SE $\beta$ Sales	Radius CI % Sales	95% Sales	Spent Per Exposed	Margin	SE ROI	H0: ROI=-100%	Ha: ROI=0%	E[t]	Mult.	E[t]=3	H0: ROI=0%	Ha: ROI=50%	E[t]	Mult.	E[t]=3	H0: ROI=0%	Ha: ROI=10%	E[t]	Mult.	E[t]=3	H0: ROI=0%	Ha: ROI=5%	E[t]	Mult.	E[t]=3
R 1	1	\$ 0.193	4.0%		\$0.16	50%	61%	1.64	3.3x	0.82	13.4x	0.16	335x	0.08	1338x	0.08	1338x	0.08	1338x	0.08	1338x	0.08	1338x	0.08	1338x	0.08	1338x	0.08	1338x
R 1	2	\$ 0.226	4.2%		\$0.06	50%	193%	0.52	33.5x	0.26	133.8x	0.05	3345x	0.03	13382x	0.03	13382x	0.03	13382x	0.03	13382x	0.03	13382x	0.03	13382x	0.03	13382x	0.03	13382x
R 1	3	\$ 0.143	5.8%		\$0.09	50%	84%	1.19	6.3x	0.60	25.2x	0.12	631x	0.06	2524x	0.06	2524x	0.06	2524x	0.06	2524x	0.06	2524x	0.06	2524x	0.06	2524x	0.06	2524x
R 1	1-6	\$ 0.912	1.4%		\$0.34	50%	134%	0.75	16.2x	0.37	64.7x	0.07	6939x	0.02	27756x	0.02	27756x	0.02	27756x	0.02	27756x	0.02	27756x	0.02	27756x	0.02	27756x	0.02	27756x
R 1	1	\$ 0.244	4.2%		\$0.04	50%	278%	0.36	69.4x	0.37	277.6x	0.07	425x	0.07	1700x	0.07	1700x	0.07	1700x	0.07	1700x	0.07	1700x	0.07	1700x	0.07	1700x	0.07	1700x
R 1	2-3	\$ 0.207	2.3%		\$0.12	50%	84%	1.20	6.3x	0.60	25.2x	0.12	629x	0.06	2515x	0.06	2515x	0.06	2515x	0.06	2515x	0.06	2515x	0.06	2515x	0.06	2515x	0.06	2515x
R 2	1a	\$ 0.139	0.9%		\$0.09	15%	24%	4.12	0.5x	2.06	2.1x	0.41	53x	0.21	212x	0.21	212x	0.21	212x	0.21	212x	0.21	212x	0.21	212x	0.21	212x	0.21	212x
R 2	1b	\$ 0.142	0.9%		\$0.08	15%	25%	3.99	0.6x	2.00	2.3x	0.40	57x	0.20	226x	0.20	226x	0.20	226x	0.20	226x	0.20	226x	0.20	226x	0.20	226x	0.20	226x
R 2	1c	\$ 0.131	0.8%		\$0.09	15%	23%	4.33	0.5x	2.17	1.9x	0.43	48x	0.22	192x	0.22	192x	0.22	192x	0.22	192x	0.22	192x	0.22	192x	0.22	192x	0.22	192x
R 3	1	\$ 0.061	9.5%		\$0.04	30%	52%	1.92	2.4x	0.96	9.7x	0.19	243x	0.10	972x	0.10	972x	0.10	972x	0.10	972x	0.10	972x	0.10	972x	0.10	972x	0.10	972x
R 3	2	\$ 0.044	8.0%		\$0.04	30%	34%	2.96	1.0x	1.48	4.1x	0.30	103x	0.15	411x	0.15	411x	0.15	411x	0.15	411x	0.15	411x	0.15	411x	0.15	411x	0.15	411x
R 3	3	\$ 0.065	6.7%		\$0.09	30%	22%	4.50	0.4x	2.25	1.8x	0.45	44x	0.23	177x	0.23	177x	0.23	177x	0.23	177x	0.23	177x	0.23	177x	0.23	177x	0.23	177x
R 3	4	\$ 0.051	7.8%		\$0.06	30%	26%	3.82	0.6x	1.91	2.5x	0.38	62x	0.19	247x	0.19	247x	0.19	247x	0.19	247x	0.19	247x	0.19	247x	0.19	247x	0.19	247x
R 3	5	\$ 0.049	5.5%		\$0.07	30%	21%	4.73	0.4x	2.36	1.6x	0.47	40x	0.24	161x	0.24	161x	0.24	161x	0.24	161x	0.24	161x	0.24	161x	0.24	161x	0.24	161x
R 3	6	\$ 0.064	4.8%		\$0.08	30%	25%	3.98	0.6x	1.99	2.3x	0.40	57x	0.20	227x	0.20	227x	0.20	227x	0.20	227x	0.20	227x	0.20	227x	0.20	227x	0.20	227x
R 3	7	\$ 0.032	10.6%		\$0.11	30%	9%	11.32	0.1x	5.66	0.3x	1.13	7x	0.57	28x	0.57	28x	0.57	28x	0.57	28x	0.57	28x	0.57	28x	0.57	28x	0.57	28x
R 4	1	\$ 0.031	10.9%		\$0.13	40%	10%	10.45	0.1x	5.22	0.3x	1.04	8x	0.52	33x	0.52	33x	0.52	33x	0.52	33x	0.52	33x	0.52	33x	0.52	33x	0.52	33x
R 5	1	\$ 0.215	0.8%		\$0.11	30%	57%	1.76	2.9x	0.88	11.6x	0.18	291x	0.09	1165x	0.09	1165x	0.09	1165x	0.09	1165x	0.09	1165x	0.09	1165x	0.09	1165x	0.09	1165x
R 5	2	\$ 0.190	4.4%		\$0.39	30%	15%	6.90	0.2x	3.45	0.8x	0.69	19x	0.34	76x	0.34	76x	0.34	76x	0.34	76x	0.34	76x	0.34	76x	0.34	76x	0.34	76x

New Accounts Only		Key Statistical Properties of Campaign										HARD			HARDER			HARDEST			CRAZY								
		Adv	#	SE $\beta$ New Accts	Radius %New Accts	95% Accts	Spent Per Person	Lifetime Value	SE ROI	H0: ROI=-100%	Ha: ROI=0%	E[t]	Mult.	E[t]=3	H0: ROI=0%	Ha: ROI=50%	E[t]	Mult.	E[t]=3	H0: ROI=0%	Ha: ROI=10%	E[t]	Mult.	E[t]=3	H0: ROI=0%	Ha: ROI=5%	E[t]	Mult.	E[t]=3
F 1	1a	69	15.6%		\$0.06	\$1,000	138%	0.73	17.1x	0.36	68.3x	0.07	1707x	0.04	6828x	0.04	6828x	0.04	6828x	0.04	6828x	0.04	6828x	0.04	6828x	0.04	6828x	0.04	6828x
F 1	1b	69	17.7%		\$0.07	\$1,000	137%	0.73	17.0x	0.36	67.9x	0.07	1697x	0.04	6790x	0.04	6790x	0.04	6790x	0.04	6790x	0.04	6790x	0.04	6790x	0.04	6790x	0.04	6790x
F 1	1c	70	10.9%		\$0.07	\$1,000	93%	1.07	7.8x	0.54	31.4x	0.11	785x	0.05	3139x	0.05	3139x	0.05	3139x	0.05	3139x	0.05	3139x	0.05	3139x	0.05	3139x	0.05	3139x
F 1	1d	70	10.5%		\$0.07	\$1,000	93%	1.08	7.7x	0.54	30.9x	0.11	774x	0.05	3094x	0.05	3094x	0.05	3094x	0.05	3094x	0.05	3094x	0.05	3094x	0.05	3094x	0.05	3094x
F 2	1	288	5.5%		\$0.03	\$1,000	47%	2.13	2.0x	1.06	8.0x	0.21	199x	0.11	795x	0.11	795x	0.11	795x	0.11	795x	0.11	795x	0.11	795x	0.11	795x	0.11	795x
F 2	1	46	8.3%		\$0.02	\$1,000	233%	0.43	48.7x	0.21	195.0x	0.04	4874x	0.02	19496x	0.02	19496x	0.02	19496x	0.02	19496x	0.02	19496x	0.02	19496x	0.02	19496x	0.02	19496x

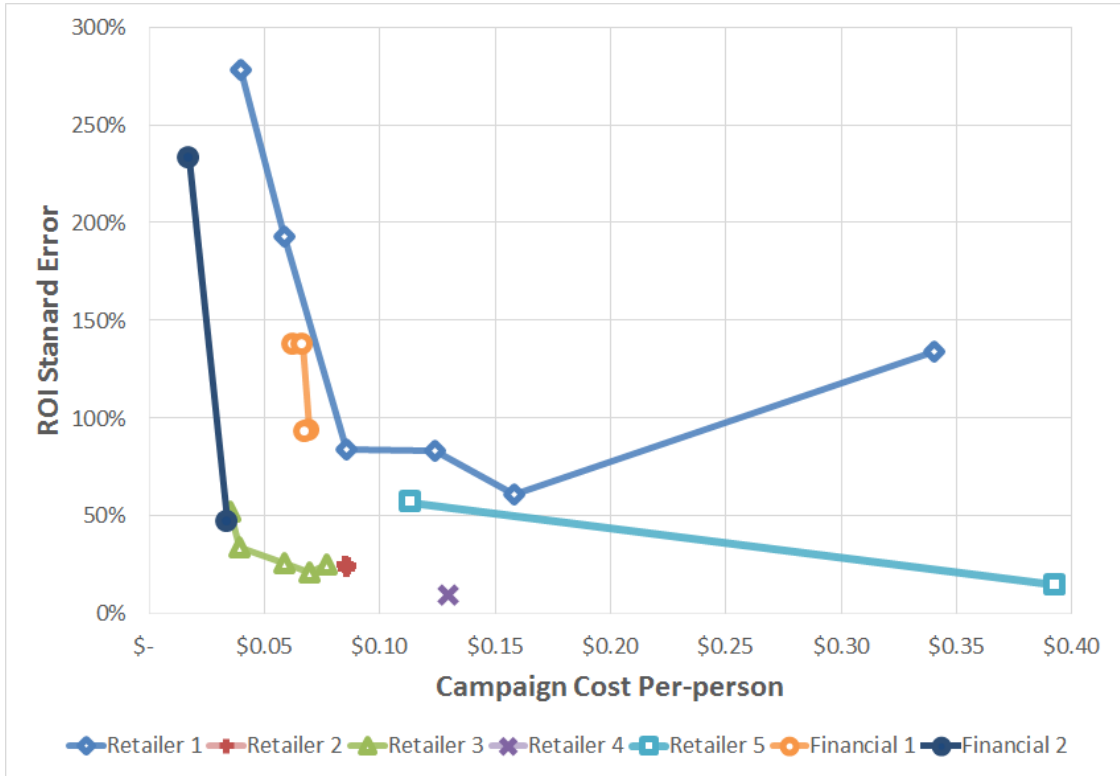


Figure 3: Relationship between ROI uncertainty and campaign cost.

what we would expect. For a given firm, a more expensive campaign requires a larger impact on sales to deliver the same percentage return. Measured against the same background noise, a larger impact is easier to identify than a smaller one—the more intense the experiment, the better the power. This implies that identifying the full shape of the  $\beta(c)$  function is effectively impossible because as one moves closer to the origin, the estimation uncertainty increases dramatically.

In the final 8 columns of Table 2, we examine an advertiser’s ability to evaluate various sets of hypotheses on the returns to expenditure. We start with disparate null and alternative hypotheses and then draw the hypotheses closer together, to tolerances more typical of investment decisions. For each hypothesis set, we give the expected  $t$ -statistic,  $E[t]$ , to reject the null hypothesis, which is a natural measure of expected statistical significance when true state of the world is given by the alternative hypothesis. An expected  $t$ -statistic of 3 provides power of 91% with a one-sided test size of 5%. We also give a “data multiplier,” which tells us how



much larger the experiment (and implicitly the total cost) would have to be in terms of new (independent) individuals to achieve  $E[t] = 3$  when the alternative hypothesis is true. The experiment could also be made larger by holding  $N$  constant and lengthening the duration using the same spend per week. Here we focus on  $N$  because it does not require us to model the within-person serial correlation of purchases. Naturally if individuals’ purchases were independent across weeks, then adding a person-week could be done just as effectively by adding another week to the existing set of targeted individuals.<sup>11</sup>

We start with distinguishing no impact (-100% ROI) from positive returns (ROI > 0%). In fact, most papers on ad effectiveness use this as the primary hypothesis of interest—the goal being to measure whether the causal influence on sales is significantly different from zero (Bagwell, 2005). Nine of 25 experiments had  $E[t] < 1.65$  (Column 9), meaning the most likely outcome was failing to reject -100% ROI when the truth was the ad was profitable.<sup>12</sup> Ten experiments had  $E[t] > 3$ , meaning they possessed sufficient power to reliably determine if the ads had a causal effect on consumer behavior. The remaining 6 experiments were moderately underpowered.

Simply rejecting that a campaign was a total waste of money is not a very ambitious goal. In the “harder” column we ask a more appropriate question from a business perspective, “Are the ads profitable?” Here we set the null hypothesis as ROI=0% and the alternative to a blockbuster return of 50%. Here 12 of 25 experiments had  $E[t] < 1$  (severely underpowered), four had  $E[t] \in [1, 2]$ , five had  $E[t] \in [2, 3]$  (90% > power > 50%), and only three had  $E[t] > 3$ . Thus, only three of the 25 had sufficient power to reliably conclude that a *wildly profitable* campaign was worth the money, and an additional 5 could reach this mark by increasing the size of the experiment by a factor of about 2.5 (those with  $E[t] \in [2, 3]$ ) or by using other methods to optimize the experimental design. The median campaign would have to be 9 times larger to have sufficient power in this setting.

The most powerful experiments were Retailer 5’s second campaign, which cost

---

<sup>11</sup>If the serial correlation is large and positive (negative), then adding more weeks is much less (more) effective than adding more people. Note also that campaigns are typically short because firms like to rotate the creative so that ads do not get stale and ignored.

<sup>12</sup>If  $E[t] < 1.65$ , even with a one-sided test, more than half the time the  $t$ -statistic will be less than the critical value due to the symmetry of the distribution. As an aside, we note that these experiments are not meant to represent optimal experimental design. Often the advertisers came to us looking to understand how much can be learned via experimentation, given a number of budgetary and campaign-objective constraints.

\$180,000 and reached 457,968 people, and Retailer 4’s campaign, which cost \$90,000 and reached 1,075,828 people. For Retailer 5’s second campaign, the relatively high precision is largely due to it having the most intense in terms of per-person spend (\$0.39). The precision improvement associated with tripling the spend as compared to an earlier campaign is shown graphically in Figure 3. Retailer 4 had good power due to two key factors: it had the fourth highest per-person spend and the second lowest standard deviation of sales.

Distinguishing a highly successful campaign from one that just broke even is not an optimization standard we typically apply in economics, yet our analysis shows that reliably distinguishing a 50% from 0% ROI is typically not possible with a \$100,000 experiment. In the third and fourth columns from the right, we draw the hypotheses closer together to a difference of 10 percentage points. While we use 0% and 10% for instructive purposes, in reality the ROI goal would need to be estimated as well (we discuss this later). Strikingly, every experiment is *severely* underpowered to reject 0% ROI in favor of 10%.  $E[t]$  is less than 0.5 for 21 of 25 campaigns, and even the most powerful experiment would have to be 7 times larger to have sufficient power to distinguish this difference. The median retail sales experiment would have to be *61 times larger* to reliably detect the difference between an investment that, using conventional standards, would be considered a strong performer (10% ROI) and one that would be not worth the time and effort (0% ROI). For new account sign-ups at financial service firms, the median multiplier is a whopping *1241*—this reflects the all-or-nothing nature of consumption patterns for these firms, a feature shared by other heavily advertised goods such as automobiles.

In the final two columns of Table 2, we push the envelope further, setting the difference between the test hypotheses to 5 percentage points. The expected  $t$ -statistics and multipliers for  $E[t] = 3$  demonstrate that this is not a question an advertiser could reasonably hope to answer for a specific campaign or in the medium-run across campaigns—in a literal sense, the total U.S. population and the advertiser’s annual advertising budget are binding constraints in most cases. These last two hypotheses sets are not straw men. These are the real standards we use in textbooks, teach our undergraduates and MBAs, and employ for many investment decisions. That they are nearly impossible to apply for these retailers and financial service providers is the key contribution of the paper, and in the discussion section we use data from industry groups to argue that these advertisers are not atypical. In fact, 5% ROI

in our setting is for a two-week period, which corresponds to an annualized ROI of over 100%. If we instead focused on 5% *annualized* ROI, the problem would be 676 times harder.<sup>13</sup>

Many investment decisions involve underlying certainty. In drug discovery, for example, a handful of drugs like *Lipitor* are big hits, and the vast majority never make it to clinical trials. Drug manufacturers typically hold large diversified portfolios of compounds for this very reason and ex-post profit measurement is relatively straightforward. Advertisers tend to vary ad copy and campaign style to diversify expenditure. And while this does guard against idiosyncratic risk of a “dud” campaign, it does not guarantee the firm is at a profitable point on the  $\beta$  function because ex-post measurement is so difficult. A good analog may be management consulting. Bloom et al. (2013) document the difficulty in measuring the returns to consulting services and conduct the first randomized trial to measure the causal influence of these expensive services. The authors report a positive effect of consulting but also report that precise ROI statements are difficult to make. One might have thought that the ability to randomize over millions of users would set advertising off to a considerable inference advantage, but this turns out not to be the case.

### 3.3 Determining the ROI target

Here we briefly touch on how a firm would determine the ROI target in our simple model. Returning to Figure 1, there are 3 important regions trifurcated by  $c^*$  and  $c_h$ .  $c^*$  gives the optimal per-person spend and defines the ROI target:  $\frac{\beta(c^*)m - c^*}{c^*}$ .  $c_h$  gives the break-even point at which average ROI is zero. For  $c < c^*$ , average ROI is positive but the firm is under-advertising—ROI is too high. For  $c > c^*$ , the firm is over-advertising, average ROI is still positive as long as  $c < c_h$ , but marginal returns are negative. In this region, although ROI is positive, spending should be reduced (which may interact with the decision maker’s average/marginal bias (de Bartolome, 1995)). When  $c > c_h$ , ROI is negative and plan of action is much clearer.

A seemingly straightforward strategy to estimate the sales impact function would be to run an experiment with several treatments in which cost per person is exogenously varied.<sup>14</sup> Each treatment gives an estimate in  $(c, \beta(c))$  space shown in Figure 1.

<sup>13</sup>We are trying to estimate 1/26th of the previous effect size, which is  $26^2$  times harder.

<sup>14</sup>See Johnson, Lewis, and Reiley (2013) for an example.

Our analysis shows that each of these points would have large confidence intervals, so fitting the function with non-parametric techniques that take into account statistical uncertainty would result in a wide range of ROI curves that cannot be rejected, providing little guidance for the advertiser.

Instead, a firm may use simple comparisons to measure marginal profit. Consider two spend levels  $0 < c_1 < c_2$ . Marginal profit is given by  $m * (\beta(c_2)m - c_2) - (\beta(c_1) - c_1) = m * (\beta(c_2) - \beta(c_1)) - (c_2 - c_1)$ . Estimating marginal ROI turns out to be more difficult primarily because the cost differential between the two campaigns  $\Delta c = c_2 - c_1$  is naturally smaller than a standalone campaign. Smaller cost differences push our effective cost per user down, meaning the points on the left side of Figure 3 are more representative of the hypothesis tests we will run with small  $\Delta c$ . To see this more clearly, notice that the variance of marginal ROI has the cost differential in the denominator:

$$Var(ROI(\Delta c)) = \left(\frac{m}{\Delta c}\right)^2 Var(\beta(c_2) - \beta(c_1)). \quad (9)$$

As we draw the two campaigns closer together to estimate marginal returns the variance of our ROI estimates diverge to infinity. This is exacerbated by the fact that the expected profit differential also decreases in  $\Delta c$  due to the concavity of  $\beta(c)$ . Ideally, we would like to find  $c^*$  where this marginal profit estimate is zero (or equal to the cost of capital), but achieving such precise estimates is essentially impossible.

In the real world various other factors further complicate matters. Concerned with ad copy “wear out,” firms tend to use different/new ad copy for different campaigns (Eastlack Jr and Rao, 1989; Lewis, 2010). Comparing campaigns of differing intensity and ad copy adds a non-trivial economic wrinkle. Using the same logic as above, determining if two creatives are significantly different will only be possible when their performance differs by a very wide margin. If creatives show considerable differences in user impact, then aggregating data across campaigns is naturally less useful, a catch-22 that places us back in the unenviable position of evaluating campaigns in isolation.

## 4 Discussion

In this section we will first address the generalizability of our findings. We then discuss implications for the broader marketplace.

### 4.1 Representativeness of our experiments

A natural concern is that our experiments are not representative of the inference challenges facing most advertisers. It is thus necessary to establish, in considerable detail, that our sample is sufficient to support the strength of our central claims. We do so now.

#### 4.1.1 Are these ads too cheap?

The higher the per capita cost of an ad, the higher the required conversion rate to break even (the relationship is linear). The ads in our study were representative of premium online display, about 1/3 the price per impression of a 30-second TV commercial, and campaigns typically delivered many impressions per user. Display ads with higher levels of targeting, or search ads, which can be viewed as highly targeted because the user has revealed a fairly specific interest, tend to be priced higher. In section 4.1.4 we discuss in detail how targeting impacts the inference problem. So thus while there are more expensive ad formats out there, these campaigns have a cost-per-user that is common in the industry and close to magnitudes of other popular advertising formats.

#### 4.1.2 Do these firms have unusually high sales volatility?

To get an idea of how the sales volatility of our firms compares to other heavily advertised categories, we use data from an industry that aggressively advertises and for which data are available: American automakers. We back out sales volatility using published data and a few back-of-the-envelope assumptions. It is reasonable to suppose the average American purchases a new car every 5–10 years. We will generously assume it is every 5 years (a higher purchase frequency makes inference easier). Suppose that the advertiser has a 15% market share. Then the annual probability of purchase for this automaker is 0.03 ( $\Pr(\text{buy}) = .2 \cdot .15 = .03$ ), which implies a standard deviation of  $\sqrt{0.03} \approx \frac{1}{6}$ . To convert this into a dollar figure,

we use the national average sales price for new cars, \$29,793.<sup>15</sup> Mean annual sales per-person is \$893 ( $\mu = \$29,793 * 0.2 * 0.15$ , price  $\times$  annual purchase rate  $\times$  market share) and  $\sigma = 1/6 * \$29,793 = \$4,700$ . This gives a  $\frac{\sigma}{\mu}$  ratio of roughly 5. However this is *yearly*, as opposed to the finer granularity used in our study. To convert this into a monthly figure, we multiply by  $(1/\sqrt{12})/(1/12) = \sqrt{12}$ , yielding a ratio of 20:1, greater than nearly all our firms.<sup>16</sup>

Heavily advertised categories such as high-end durable goods, subscription services such as credit cards, and infrequent big-ticket purchases like vacations all seem to have consumption patterns that are more volatile than the retailers we studied selling sweaters and dress shirts and about as volatile as the financial service firms who also face an “all-or-nothing” consumption profile. Political advertising appears to share similar difficulties (Broockman and Green, 2013).<sup>17</sup>

It is important to note that our results do not necessarily apply to small firms, brand new products or direct-response TV advertising. However, according to estimates from Kantar AdSpender (and other industry sources), large advertisers using standard ad formats, such as the ones we study, account for the vast majority of advertising expenditure. Thus while our results do not apply to every market participant, they do have important implications for the market generally.

### 4.1.3 Are these campaigns too small?

Our scale multipliers give an idea of the cost necessary to push confidence intervals to informative widths—the implied cost (if that many unique individuals were avail-

<sup>15</sup>Source: <http://www.nada.org/Publications/NADADATA/2011/default>.

<sup>16</sup>Here we assumed zero variance in purchase price. In this setting, including variation in price does not make the inference problem much more difficult—most of the difficulty is driven by rarity of purchases. The variance of each component contributes to overall sales volatility as given by: Let  $Y = p * \$$  where  $p$  is purchase probability and  $\$$  is basket size.  $\frac{\sqrt{\text{var}(Y)}}{E[Y]} = \frac{\sqrt{E[p]^2 \text{Var}(\$) + \text{Var}(p) E[\$]^2 + \text{Var}(\$) \text{Var}(p)}}{E[p] E[\$]}$ . Both components can presumably be impacted by advertising. For the retailers in our study, back-of-the-envelope calculations indicate that each component contributes significantly to the total, but purchase rarity probably accounts for a larger portion than variation in basket size conditional on purchasing. For example, using values from Lewis and Reiley (2013), calculations show that ignoring the components with  $\text{Var}(\$)$  reduces the total coefficient of variation by about 40%. Ignoring the basket size component would make the problem somewhat easier statistically, but would induce a bias of unknown size.

<sup>17</sup>Their figures imply that one would need 400,000 unique users, to reliably reject a cost of \$50 per marginal vote if the ads, in fact, have no effect. Even this coarse test would not be feasible for many candidates in many elections.

able) was often in the tens of millions of dollars, far more expensive than even the largest reach advertisement in the US, the NFL Super Bowl, which we will use here in a thought experiment, supposing that the 30-second TV spots can be individually randomized. We will try to define the set of advertisers that can both afford a Super Bowl spot and detect the return on investment.

The affordability constraint is simply an accounting exercise to ensure firm’s advertising budget can accommodate such a large expenditure. For the budget we choose a value, 5% of revenue, which exceeds advertising budgets for most major firms.<sup>18</sup> To build intuition on the detectability constraint, recall that ROI is the percentage return on the *ad cost*—it does not depend on the baseline level of sales. The sales *level* lift that nets a positive ROI is a much larger *percentage* lift for a small firm than for larger firms and thus more likely to stand out statistically. The “detectability constraint” gives the largest firm, in terms of annual revenue, that can meaningfully evaluate a given ROI hypothesis set.

Out of consideration for space, we put the formal argument in the Appendix. We set the analysis window  $w = 2$  (weeks) to match most of the analysis of this paper,  $t_{ROI} = 3$  to match our standard power requirement, and  $\frac{\sigma}{\mu} = 10$  to match the value we see strong evidence for in our study, even though it will understate volatility for advertisers such as automakers and financial service firms.<sup>19</sup> We report bounds for two values of gross margin: 0.25 and 0.50. The final step is to calibrate pricing and audience. We use the following parameters:  $N_E$  is 50 million (1/2 the viewers) and the cost of the ad is 1/2 the market rate,  $C = \$1,000,000$  (1/2 the cost).

Table 3 gives the upper and lower bounds on annual revenue. If an ad promotes only a specific product group, for instance the 2011 Honda Civic, then the relevant figure to compare to the bounds would be the revenue for that product group. Examining Row 1, we see that most companies would be able to reliably determine if the ad causally impacted consumers. Major automobile manufacturers (which are low margin) doing brand advertising would exceed this limit, but specific model-years fall below it.<sup>20</sup>

We see in Row 2 that many companies and product categories could reliably

---

<sup>18</sup>Source: Kantar AdSpender.

<sup>19</sup>We use  $\rho = 0.5$  to match the empirical viewing share for adults for the Super Bowl. See Appendix for more details.

<sup>20</sup>However, we have assumed a  $\frac{\sigma}{\mu}$  ratio of 10, which is probably half the true value for car sales over 2-4 week time frame, meaning the correct bound is probably twice as high.

Table 3: Super Bowl “Impossibility” Theorem Bounds

$H_A$ : ROI	$H_0$ : ROI	Affordability Annual Rev.	Detectability, $m=.50$ Annual Rev.	Detectability, $m=.25$ Annual Rev.
0%	-100%	\$2.08B	\$34.47B	\$63.3B
50%	0%	\$2.08B	\$17.33B	\$34.6B
10%	0%	\$2.08B	\$3.47B	\$6.9B
5%	0%	\$2.08B	\$1.73B	\$3.4B

distinguish 50% ROI from 0%—the bounds are \$17.3 billion and \$34.6 billion for the high and low margins respectively—but large firms or products could not. For the final two hypothesis sets, the bands are tight to vanishing. It is nearly impossible to be large enough to afford the ad, but small enough to reliably detect meaningful differences in ROI.

#### 4.1.4 Are these campaigns not targeted enough?

Can a firm more powerfully assess their advertising stock by performing experiments on the particularly susceptible portion of the population? Suppose there are  $N$  individuals in the population the firm would consider advertising to. We assume that the firm does not know how a campaign will impact each individual, but can order them by expected impact. The firm wants to design an experiment using the first  $M$  of the possible  $N$  individuals. We define  $\Delta\mu(M)$ ,  $\sigma(M)$ , and  $c(M)$  as the mean sales impact, standard deviation of sales, and average cost functions, respectively, when advertising to the first  $M$  people. The  $t$ -statistic against the null hypothesis of -100% ROI is given by:  $t = \sqrt{\frac{M}{2}} \cdot \frac{\Delta\mu(M)}{\sigma(M)}$ .

Assuming constant variance,  $\sigma^2(M) = \sigma^2$ , and taking the derivative with respect to  $M$ , we get:

$$\frac{dt}{dM} = \frac{1}{2\sqrt{2M}} \frac{\Delta\mu(M)}{\sigma} + \sqrt{\frac{M}{2}} \frac{\Delta\mu'(M)}{\sigma} \quad (10)$$

With targeting,  $\Delta\mu'(M) < 0$ . Simplifying the right hand side, we find the  $t$ -statistic is increasing in  $M$  if the targeting effect decays slower than  $\frac{\Delta\mu(M)}{2\sqrt{2M}}$ . Thus, the question of whether targeting helps or hurts inference is an empirical one. If the sales impact is concentrated on a certain portion of the population, one is better off reducing sample size to gain a higher signal-to-noise ratio. Conversely, if influence



is spread rather evenly across the population, targeting damages power. Additional details of this argument are in the Appendix.

#### 4.1.5 Would longer measurement windows help?

Any analysis of the returns to advertising invariably has to specify the window of time to be included in the study. We followed the standard practice of the campaign period and a relatively short window after the campaign ended. Perhaps by adding more data on the time dimension, we would get a better estimate of the cumulative impact and improve statistical precision.

We present the formal argument in the Appendix. The key proposition is the following:

*If the next week's expected effect is less than one-half the average effect over all previous weeks, then adding it in will only reduce the t-statistic.*

The proposition tells us when a marginal week hurts estimation precision because it introduces more noise than signal. As an example, suppose the causal impact of the advertising on weeks 1, 2, and 3 is 5%, 2%, and  $z$ , respectively. Then  $z$  must be greater than  $\frac{5+2}{2} = 1.75$ . In other words, unless there is very limited decay in the ad effect over time, we would be better off curtailing the evaluation window to two weeks. With moderate decay, optimal evaluation windows (from a power perspective) get quite short. An additional week of data increases the effective sample size and the cumulative impact, but reduces the average per-time-period impact, watering down the effect we are trying to measure. The proposition can provide helpful guidance and helps explain why short windows are generally used, but quantitatively applying it requires precise ROI estimates for the very inference problem we are trying to solve.

#### 4.1.6 What about sponsored search advertising?

Sponsored search is the practice of paying to place links at the top of a search engine results page. Given these ads are highly targeted, they typically cost more per impression than display ads. All else equal, the high per capita costs makes the inference problem easier. However all else is not equal. First, one cannot arbitrarily increase sample size because it is driven by the frequency of the chosen query. This

places the problem in our targeting framework already discussed. Second, large advertisers are typically relevant to the query and thus show up somewhere in the “organic results.” For instance, for the query “car insurance” one will immediately notice the advertisers in the sponsored links—major insurance firms—are very high in the organic web results directly below the ads. The fact that an advertiser gets both organic and paid clicks makes it hard to determine how many of the paid clicks are incremental. Note that this problem would not be present for an advertiser that does not have high organic relevance. The inference challenge is exacerbated considerably by the fact that search engines encrypt organic clicks, which means that while one is able to tell a click came from google.com, for instance, she could not to tell what the user had searched for.<sup>21</sup> This severely limits the type of experimentation that is possible and in particular, eliminates user-based randomization (because it is not possible to form the unexposed group).

Given these limitations, experiments typically randomize over temporal or geographic units. Recent work by eBay Research Labs (Blake, Nosko, and Tadelis, 2013) conducts a very large experiment—easily in the tens of millions of dollars—using these methods.<sup>22</sup> The authors examine the returns to branded keywords (e.g., “tablet computer *ebay*”) and unbranded keywords. For branded terms, pause experiments show that most of the clicks on paid search links would have otherwise occurred on an organic link. For unbranded terms, geo-randomization is used to estimate that paid search is causally linked to 0.44% of total sales, with a standard error of 0.62%, leading to a 95% confidence interval of (-0.77%, 1.66%). The 0.44% sales impact corresponds to an *average* ROI of -68%, a considerable loss; however, the top of the confidence interval is +16% ROI, meaning they could not reject profitability at standard confidence levels. These confidence intervals are similar to our median experiment, which highlights the importance of individual-based randomization.

To understand the impact of a *marginal* dollar, the authors regress sales revenue on search spending using the randomization as the instrument. Ordinary Least Squares, even with a full set of controls, grossly overstates the true impact due to temporally varying purchase intent, a bias first documented in Lewis et al. (2011).

---

<sup>21</sup>Encryption removes the query string and other parameters from the “referring URL” but leaves the “top level domain.”

<sup>22</sup>As is true of all the content in this paper, the views expressed are solely our own and not those of our employers, both of which operate search engines.

The Instrumental Variables estimate had a 95% confidence interval about 30 times wider than the point estimate.

We think this paper dovetails our results nicely. First, the authors are employed by a large advertiser and openly claim the company did not know the returns to advertising and strongly imply (p. 14, paragraph 2) that observational methods were being used that severely overstated returns. Second, the experiment confirms that truly ineffective campaigns can be identified via large scale experimentation (but not observational methods). Third, the estimates on the marginal dollar spent have enormous confidence intervals, and the considerably smaller confidence interval on the average ROI is still over 100 percentage points wide.

## 4.2 Improving experimental design and data collection

### 4.2.1 Optimizing ex-post evaluation and experimental design

In Section 2.2 we calibrated an advertiser’s inference problem using univariate linear regression for expository clarity. In our actual estimation, we conditioned on the user level covariates listed in the column labeled by the vector  $\mathbf{W}$  in Table 1 using several methods to strengthen power; such panel techniques predict and absorb residual variation. Lagged sales are the best predictor and are used wherever possible, reducing variance in the dependent variable by as much as 40%. This is echoed by a recent paper on improving the power of online experiments, Deng et al. (2013), which finds that lagged dependent variables can reduce residual variation by as much as 50%. However, seemingly large improvements in  $R^2$  lead to only modest reductions in standard errors. A little math shows that going from  $R^2 = 0$  in the univariate regression to  $R^2_{|\mathbf{W}} = 50\%$  yields a sublinear reduction in standard errors of 29%. Hence, the modeling is as valuable as doubling the sample—a significant improvement, but one that does not materially change the measurement difficulty.<sup>23</sup>

Another method to improve power is to use concentrated tests. Figure 3 shows that larger expenditures per-person (conditional on experiment size) are associated

---

<sup>23</sup> $1 - \sqrt{\frac{1-R^2_{|\mathbf{W}}}{1-R^2}} = 1 - \sqrt{1 - R^2_{|\mathbf{W}}} = 29\%$ . An order-of-magnitude reduction in standard errors would require  $R^2_{|\mathbf{W}} = 99\%$ , perhaps a “nearly impossible” goal. From a design perspective, a related method is to match like users into pairs and then randomize exposure within each pair. This sort of pre-experiment matching is more useful for small samples as it can insure against an unlucky randomization step—with large samples this is highly unlikely to occur. Demographics are typically not found to meaningfully improve power.

with lower standard errors on the ROI estimate. We previously showed that a firm can sometimes improve power by using the most susceptible part of the population in the experiment. Similarly, if a portion of the customer base had lower variance in sales over time, these customers are attractive to include in an experiment. In these ways, firms can carve out a portion of the population with statistically favorable properties. Of course, the measurements of the returns to advertising for these subsamples may not reflect returns for the population of consumers the firm wishes to advertise to. Moreover, a concentrated test (high per-person spend) may be at a level of advertising beyond the firm’s optimum.

The final design advance on the horizon is free control ads for experiments, sometimes referred to as “ghost ads.” In the past, advertisers would have to pay for both treatment and control ads (often an ad for a charity), so that they could measure who actually saw a control ad to form the comparison group. Technology now exists to hold the control group out from the advertiser’s ad but measure when an attempt to serve the treatment ad was made. The system records what ad was actually shown, so one can control for possible competitive effects of advertising, but this does add a non-trivial wrinkle to the analysis. Even so, the prospect of free control ads means doubling experiment size without any additional cost is feasible. Examining the multipliers in Table 2, we see this will help this class of advertisers running \$100,000 experiments infer if ad expenditure causally impacts consumers. An improvement of 2x will in general not be much help evaluating any of the more realistic hypotheses sets, but any reduction in experimental costs will only promote accurate inference in the marketplace.

Our conclusion that while the methods discussed in this section offer improvements in power for some firms in some situations, they do not solve the fundamental inference difficulties we have raised in this paper.

#### **4.2.2 Aggregating across experiments and forming priors**

If a firm was committed to evaluating advertising spend, then for the media in which experimentation is available (currently a minority of total advertising expenditure), our results indicate that running repeated \$500,000 experiments would allow some firms to understand the *average* impact of global spend. One strategy would be to use an evolving prior to evaluate campaigns. But as we have seen, the signal

from any given campaign is relatively weak, meaning a Bayesian update would essentially return the prior. So while this is a promising strategy to determine the global average, it probably would not help much in evaluating campaigns. It also needs to be stressed that such a strategy requires an enormous commitment by the firm, one that does not appear to be commonplace today, and for many large advertisers even this sort of commitment would not be enough.

## 4.3 Marketplace implications

### 4.3.1 A new competitive advantage of scale

An implication of the low power of advertising experiments is that large publishers have an advantage not only through the common notion of having larger reach, but also by having the user base to run reliable experiments. Table 2 shows that many large advertisers could narrow confidence intervals to an acceptable tolerance with experiments in the tens of millions of users in each treatment cell. Only the largest publishers could offer such a product. If experimentation becomes more common, a trend we believe is occurring, then scale will increasingly confer a new competitive advantage. A smaller publisher, such as the *New York Times*, simply cannot provide same quality feedback as a massive publisher and may be better off outsourcing ad-serving to a larger network (which could then include the inventory as part of a larger experiment). For smaller advertisers, the large publisher can leverage its scale to recommend ad features based on findings from past experimentation with larger firms. Increased experimentation thus has the potential to fundamentally shape the organization of web publishing and other advertising-based industries.

### 4.3.2 The impact of noisy signals on within-firm communication

The uncertainty surrounding ROI estimates can interact with incentives to create a moral hazards in communication. Suppose the “media buyer” gets a bonus based his manager’s posterior belief on campaign ROI. If reports are delivered with certainty and completely verifiable, there is no agency problem. If they are totally unverifiable, we are in a *cheap talk game* (Crawford and Sobel, 1982) where strategic communication leads to reports that are correlated with the agent’s signal (the estimate), but noisy due to the common knowledge of the agent’s bias. Since it is very

hard to disprove a report with other data and estimates themselves are noisy and likely manipulatable<sup>24</sup> a cheap talk game might be a useful modeling approximation.

Alternatively we might view the strategic communication as a *persuasion game* (Milgrom and Roberts, 1986). Applying the model of Shin (1994), we suppose that the manager is unsure which campaigns have verifiable ROI estimates. In equilibrium, the manager will be *skeptical* because the media buyer will report good news when available but filter bad news, which limits the amount of information that can travel of up the chain of command.

### 4.3.3 Variance in advertising spend across competitors

Information is scarce in the advertising market, meaning that the “selective pressure” on advertising spending is weak. We would thus expect significant heterogeneity in advertising spend by similar firms in the same industry. Empirically testing this prediction is difficult because there are many economic reasons firms could have different advertising strategies.<sup>25</sup> We thus limit our comparisons to industries dominated by a handful of firms that share key characteristics reported to the SEC such as margins, access to technology, annual revenue, and customer base. Our data on advertising expenditure comes from Kantar Media’s AdSpender report.

In Appendix Table 1 we give advertising expenditure, revenue, and margin for the following U.S. industries: rental cars, mobile phone carriers, international airlines, online financial services, and fast food. These constitute the markets that met the requirements we have laid out and had data availability. The data reveal distinct high/low advertising strategies. Advertising expenditure as a percent of revenue differs by more than a factor of five. For example, online brokerages Scottrade, TD Ameritrade, and ETrade have similar business models and report identical gross margins. ETrade pursues a high advertising strategy, with 12.63% of revenue going to advertising. Scottrade spends 8.45% and TD Ameritrade pursues a

---

<sup>24</sup>This can be done by varying the estimation technique, changing the control variable set to find the highest point estimate, etc. With fragile point estimates these techniques can be quite “effective.” If the principal could access the raw data at some cost it could mitigate this problem. As a practical matter, it is unclear who has the incentive and time to do this. If there are multiple biased agents, all that matters is the most biased agent, who forms a choke point (Ambrus et al., 2013).

<sup>25</sup>For example, low-cost retailers might compete primarily on price and advertise very little because it erodes slim margins. As we saw in section 2, the lower a firm’s margin, the higher the impact of ad has to be to break even.

low-advertising strategy, *6.93 times less* than ETrade per dollar of revenue. We observe this pattern in most of the qualifying industries. This evidence is by no means conclusive, but the existence of vastly different advertising strategies by seemingly similar firms operating in the same market with similar margins is consistent with our prediction that vastly different beliefs on the efficacy of advertising are allowed to persist in the market.

#### 4.3.4 How unusual is this market?

In markets with limited informational feedback as to the efficacy of the product, sellers may have a customer base that holds fundamentally incorrect beliefs. Here we look at two industries that we think share this feature. The first is management consulting. Bloom et al. (2013) argue that consulting expenditures are rarely implemented in a way that the relevant counterfactual can be formed. To overcome this endogeneity problem, the authors ran a controlled experiment and documented a positive impact of consulting services, but also documented that making precise ROI statements is incredibly difficult.

The second is the vitamin and supplement market. The industry grosses about \$20 billion annually, yet it is a contentious point in the medical community as to whether supplements do *anything* for a healthy individual (the main customer base).<sup>26</sup> The Physicians Health Study II (Lee et al., 2005) followed 39,876 healthy women over 12 years. Half received vitamin E through a supplement; the other half took a placebo. The 95% confidence interval on the impact on heart attacks ranged from a 23% risk reduction to an 18% risk increase. We can translate this uncertainty into an “economic confidence interval” using a recent estimate placing cost of a heart attack around \$1 million (Shaw et al., 2006). The economic confidence interval is \$192 million wide—a whopping *100 times* the \$2.1 million cost of vitamins for the study. The economic confidence interval for cancer was of a similar magnitude.

---

<sup>26</sup>Supplements are supposed to improve health for a *healthy person*, not prevent vitamin deficiency diseases such as rickets and scurvy, because in the developed world one gets enough of these vitamins through even the unhealthiest of diets (Ward et al., 2007).

## 5 Conclusion

Using one of the largest collection of field experiments to-date, we have shown that inferring the effects of advertising is exceedingly difficult. These findings have deep industrial organization implications. First, the advertising market as a whole may have incorrect beliefs about the causal impact of advertising on consumer behavior. As experimentation becomes more common and some firms commit the resources to run the massive (or many large, repeated) experiments necessary to generate informative signals, there could be a meaningful shift in advertising prices. Second, weak signals mean priors can dominate decision making, helping to explain why advertising spending varies widely across similar firms in the same industry. Third, the requirement for huge sample sizes in experimentation sets the largest publishers off to an advantage—if the market begins to demand information, their scale will pay an “informational dividend.” Overall, this data landscape means advertisers’ decision-making differs from our standards notion of profit maximization, fundamentally shaping the market for advertising.

## References

- Abraham, M. (2008). The off-line impact of online ads. *Harvard Business Review*, 86(4):28.
- Abraham, M. and Lodish, L. (1990). Getting the most out of advertising and promotion. *Harvard Business Review*, 68(3):50.
- Ambrus, A., Azevedo, E. M., and Kamada, Y. (2013). Hierarchical cheap talk. *Theoretical Economics*, 8(1):233–261.
- Bagwell, K. (2005). The economic analysis of advertising. *Handbook of Industrial Organization Volume 3*.
- Blake, T., Nosko, C., and Tadelis, S. (2013). Consumer heterogeneity and paid search effectiveness: A large scale field experiment. *NBER Working Paper*, pages 1–26.
- Bloom, N., Eifert, B., Mahajan, A., McKenzie, D., and Roberts, J. (2013). Does management matter? Evidence from India. *The Quarterly Journal of Economics*, 128(1):1–51.



- Broockman, D. E. and Green, D. P. (2013). Do online advertisements increase political candidates' name recognition or favorability? Evidence from randomized field experiments. *Political Behavior*.
- Card, D. (1999). The causal effect of education on earnings. *Handbook of Labor Economics*, 3:1801–1863.
- Carroll, V., Rao, A., Lee, H., Shapiro, A., and Bayus, B. (1985). The Navy enlistment marketing experiment. *Marketing Science*, 4(4):352–374.
- Crawford, V. P. and Sobel, J. (1982). Strategic information transmission. *Econometrica: Journal of the Econometric Society*, pages 1431–1451.
- de Bartolome, C. A. (1995). Which tax rate do people use: Average or marginal? *Journal of Public Economics*, 56(1):79–96.
- Deng, A., Xu, Y., Kohavi, R., and Walker, T. (2013). Improving the sensitivity of online controlled experiments by utilizing pre-experiment data. In *Proceedings of the sixth ACM international conference on Web search and data mining*, pages 123–132. ACM.
- Eastlack Jr, J. and Rao, A. (1989). Advertising experiments at the Campbell Soup Company. *Marketing Science*, pages 57–71.
- Fulgoni, G. and Morn, M. (2008). How online advertising works: Whither the click. *Comscore.com Whitepaper*.
- Gelman, A. and Carlin, J. (2013). Beyond power calculations to a broader design analysis, prospective or retrospective, using external information. *Working Paper*.
- Johnson, G., Lewis, R., and Reiley, D. (2013). Add more ads? experimentally measuring incremental purchases due to increased frequency of online display advertising. *Working paper*.
- Kaiser, H. (2005). *Economics of Commodity Promotion Programs: Lessons from California*. Peter Lang Publishing.
- LaLonde, R. (1986). Evaluating the econometric evaluations of training programs with experimental data. *The American Economic Review*, 76(4):604–620.
- Lee, I.-M., Cook, N. R., Gaziano, J. M., Gordon, D., Ridker, P. M., Manson, J. E., Hennekens, C. H., and Buring, J. E. (2005). Vitamin E in the primary

- prevention of cardiovascular disease and cancer. *The Journal of the American Medical Association*, 294(1):56.
- Lewis, R. A. (2010). *Where's the "Wear-Out?": Online Display Ads and the Impact of Frequency*. PhD thesis, MIT PhD Dissertation.
- Lewis, R. A., Rao, J. M., and Reiley, D. H. (2011). Here, there, and everywhere: correlated online behaviors can lead to overestimates of the effects of advertising. In *Proceedings of the 20th international conference on World wide web*, pages 157–166. ACM.
- Lewis, R. A., Rao, J. M., and Reiley, D. H. (Forthcoming). chapter Measuring the Effects of Advertising: The Digital Frontier. NBER Press.
- Lewis, R. A. and Reiley, D. H. (2013). Online advertising and offline sales: Measuring the effects of retail advertising via a controlled experiment on yahoo! *Unpublished manuscript*.
- Lippman, S. M., Klein, E. A., Goodman, P. J., Lucia, M. S., Thompson, I. M., Ford, L. G., Parnes, H. L., Minasian, L. M., Gaziano, J. M., Hartline, J. A., et al. (2009). Effect of selenium and vitamin E on risk of prostate cancer and other cancers. *The Journal of the American Medical Association*, 301(1):39–51.
- Lodish, L., Abraham, M., Kalmenson, S., Livelsberger, J., Lubetkin, B., Richardson, B., and Stevens, M. (1995). How TV advertising works: A meta-analysis of 389 real world split cable TV advertising experiments. *Journal of Marketing Research*, 32(2):125–139.
- Lovell, M. (2008). A simple proof of the FWL theorem. *The Journal of Economic Education*, 39(1):88–91.
- Milgrom, P. and Roberts, J. (1986). Relying on the information of interested parties. *The RAND Journal of Economics*, pages 18–32.
- Shaw, L. J., Merz, C. N. B., Pepine, C. J., Reis, S. E., Bittner, V., Kip, K. E., Kelsey, S. F., Olson, M., Johnson, B. D., Mankad, S., et al. (2006). The WISE study: The economic burden of angina in women with suspected ischemic heart disease results from the NIH National Heart, Lung, and Blood Institute sponsored womens ischemia syndrome evaluation. *Circulation*, 114(9):894–904.

Shin, H. S. (1994). News management and the value of firms. *The RAND Journal of Economics*, pages 58–71.

Ward, L. M., Gaboury, I., Ladhani, M., and Zlotkin, S. (2007). Vitamin D-deficiency rickets among children in Canada. *CMAJ*, 177(2):161–166.

Wilbur, K. (2008). How the digital video recorder (DVR) changes traditional television advertising. *Journal of Advertising*, 37(1):143–149.

## 6 Appendix

### 6.1 Super Bowl Impossibility Theorem

We will now present the formal argument and calibrate it with data from our experiments and publicly available information on Super Bowl advertising. We need to define some terms. Let  $N_{Total}$  be the total adult population,  $N$  be the total adult audience, and  $\rho = \frac{N}{N_{Total}}$  be the reach of the Super Bowl.  $N_E$  gives the number of reached (exposed) individuals; we set  $N_E = N/2$  to maximize power. On the cost side,  $C$  is the total cost of the ad, and  $c$  is the cost per exposed person. Let  $\mu$  equal the mean purchase amount for all customers during the campaign window and  $\sigma$  be the standard deviation of purchases for customers during the campaign window. We will use  $\frac{\sigma}{\mu}$ , the coefficient of variation, which we have noted is typically 10 for advertisers in our sample and greater than 10 in other industries, to calibrate the argument.  $m$  is the gross margin for the advertiser’s business

We also need to define a few terms to describe the advertiser’s budget. Let  $w$  be the number of weeks covered by the campaign’s analysis (and the advertising expense),  $b$  give the fraction of revenue devoted to advertising (% advertising budget), and  $R$  be the total annual revenue. To get the affordability bound, we define  $\gamma_C$  as the fraction of the ad budget in the campaign window devoted to the Super Bowl ad. For instance, if  $\gamma_C = 1$ , this means the firm spends all advertising dollars for the period in question on the Super Bowl.

We now present the argument, which is an algebraic exercise with one key step: substituting for the coefficient of variation and solving for the revenue bounds.

First we construct the affordability bound. To afford the ad, it must be the case that it costs less than the ad budget, which is the revenue for the time period in

question,  $R \cdot \frac{w}{52}$ , times  $b$ , the percentage of the revenue devoted to advertising, times  $\gamma_c$ , the fraction of the budget that can be devoted to one media outlet:

$$C \leq \left( R \cdot \frac{w}{52} \right) \cdot b \cdot \gamma_c.$$

Solving this equation for revenue gives the affordability limit:

$$R \geq \frac{C}{\gamma_c b \cdot \frac{w}{52}}. \quad (11)$$

For the detectability limit, let  $r$  and  $r_0$  be the target ROI and null hypothesis ROI, respectively. The  $t$ -statistic is given by:

$$\begin{aligned} t_{ROI} &\leq \frac{r - r_0}{\sqrt{\frac{2}{N}} \times \sigma_{ROI}} \\ t_{ROI} &\leq \frac{(r - r_0)}{\sqrt{\frac{2}{N}} \left( \frac{m\sigma}{c} \right)} \\ t_{ROI} &\leq \frac{(r - r_0)}{\sqrt{\frac{2}{N}} \left( \frac{\sigma}{\mu} \right) / \frac{c}{m\mu}}. \end{aligned}$$

The first equation is just the definition of the test statistic. The second equation follows from substituting in the standard deviation of ROI, which is a linear function of the sales standard deviation, per capita cost, and gross margin. The final equation simply multiplies the denominator by  $\frac{\mu}{\mu}$ . We do this so we can substitute in a constant for the coefficient of variation,  $\frac{\sigma}{\mu}$ , and solve for  $\mu$ , as given below:

$$\mu \leq \frac{(r - r_0) c}{\sqrt{\frac{2}{N}} \left( \frac{\sigma}{\mu} \right) m \cdot t_{ROI}} \equiv \bar{\mu}$$

The right-most definition is for notational convenience. We can also relate mean sales during the campaign period to total revenue:

$$\mu = R \cdot \frac{w}{N_{Total}}. \quad (12)$$

We then solve for revenue and substitute in  $\bar{\mu}$  for  $\mu$  to get the detectability limit:

$$R \leq \frac{N_{Total} \cdot \bar{\mu}}{\frac{w}{52}} \quad (13)$$

Examining the detectability limit, referring back to  $\bar{\mu}$  where necessary, we see that it decreases with  $\frac{\sigma}{\mu}$ . This is intuitive, as the noise to signal ratio increases, inference becomes more difficult. It also falls with the required  $t$  and gross margin. To understand why the bound rises as margin falls, consider two companies, one with a high margin, one with a low margin. All else equal, the low margin firm is experiencing a larger change in sales for a given ROI change. Naturally the bound also rises with the gap between the null hypothesis and target ROI.

Putting both limits together, we obtain the interval for detectability and affordability in terms of the firm's annual revenue:

$$\frac{C}{\gamma_C b \cdot \frac{w}{52}} \leq R \leq \frac{N_{Total} \cdot \bar{\mu}}{\frac{w}{52}}. \quad (14)$$

## 6.2 Targeting details

The standard deviation of the ROI,  $\sigma_{ROI}$ , is given by:

$$\begin{aligned} ROI &= \frac{\Delta\mu(M)}{C(M)} - 1 \\ \sigma_{ROI}^2 &= Var\left(\frac{\Delta\mu(M)}{C(M)}\right) = \frac{2\sigma^2(M)}{M \cdot (C(M))^2} \end{aligned}$$

which implies:

$$\sigma_{ROI} = \frac{\sigma(M)}{\sqrt{M/2} \cdot C(M)} \quad (15)$$

Notice that this formula does not rely upon the actual impact of the ads, except that we calibrate the expected effect against the cost (in reality, costs will be correlated with ad impact). It only incorporates the average volatility of the  $M$  observations. The standard error of our estimate of the ROI is decreasing in  $M$  as long as the ratio  $\sigma(M)/C(M)$  does not increase faster than  $\sqrt{M}$ . For the special case of a constant variance, the standard error of the ROI can be more precisely estimated as long as the average costs do not decline faster than  $\frac{1}{\sqrt{M}}$ . Note average

costs cannot decline faster than  $\frac{1}{M}$  unless the advertiser is actually paid to take extra impressions, which seems unlikely. Another special case is constant average cost. Here as long as  $\sigma(M)$  does not increase faster than  $\sqrt{M}$ , more precision is gained by expanding reach.

### 6.3 Campaign window proof

Note this entire argument is also in a forthcoming NBER book chapter.

We again employ the  $t$ -statistic, but also index little  $t$  for time. For the sake of concreteness, let time be indexed in terms of weeks. For notational simplicity, we will assume constant variance in the outcome over time, no covariance in outcomes over time,<sup>27</sup> constant variance across exposed and unexposed groups, and balanced group sizes. We will consider the long-term effects by examining a cumulative  $t$ -statistic (against the null of no effect) for  $T$  weeks rather than a separate statistic for each week. We write the cumulative  $t$ -statistic for  $T$  weeks as:

$$t_{\Delta\bar{y}_T} = \sqrt{\frac{N}{2}} \left( \frac{\sum_{t=1}^T \Delta\bar{y}_t}{\sqrt{T}\hat{\sigma}} \right). \quad (16)$$

At first glance, this  $t$ -statistic appears to be a typical  $O(\sqrt{T})$  asymptotic rate with the numerator being a sum over  $T$  ad effects and the denominator growing at a  $\sqrt{T}$  rate. This is where economics comes to bear. Since  $\Delta\bar{y}_t$  represents the impact of a given advertising campaign during and following the campaign (since  $t = 1$  indexes the first week of the campaign),  $\Delta\bar{y}_t \geq 0$ . But the effect of the ad each week cannot be a constant—if it were, the effect of the campaign would be infinite. Thus it is generally modeled to be decreasing over time.

With a decreasing ad effect, we should still be able to use all of the extra data we gather following the campaign to obtain more statistically significant effects, right? Wrong. Consider the condition necessary for an additional week to increase the

---

<sup>27</sup>This assumption is clearly false: individual heterogeneity and habitual purchase behavior result in serial correlation in purchasing behavior. However, as we are considering the analysis over time, if we assume a panel structure with fixed effect or other residual-variance absorbing techniques to account for the source of this heterogeneity, this assumption should not be a first-order concern.

$t$ -statistic:

$$t_{\Delta\bar{y}_T} < t_{\Delta\bar{y}_{T+1}}$$

$$\frac{\sum_{t=1}^T \Delta\bar{y}_t}{\sqrt{T}} < \frac{\sum_{t=1}^{T+1} \Delta\bar{y}_t}{\sqrt{T+1}}$$

Some additional algebra leads us to

$$1 + \frac{1}{T} < \left( 1 + \frac{\Delta\bar{y}_{T+1}}{\sum_{t=1}^T \Delta\bar{y}_t} \right)^2$$

which approximately implies

$$\frac{1}{2} \cdot \frac{1}{T} \sum_{t=1}^T \Delta\bar{y}_t < \Delta\bar{y}_{T+1}. \quad (17)$$

This last expression says, “If the next week’s expected effect is less than one-half the average effect over all previous weeks, then adding it in will only reduce precision.” Thus, the marginal week can actually cloud the previous weeks, as its signal-to-noise ratio is not sufficiently large enough to warrant its inclusion.<sup>28</sup> If the expected impact of the campaign following exposure decays rapidly (although not necessarily all the way to zero), it is likely that including additional weeks beyond the campaign weeks will decrease the statistical precision.

Suppose that you were just content with the lower bound of the confidence interval increasing in expectation. A similar calculation, under similar assumptions, shows that the lower bound of a 95% confidence interval will increase if and only if

$$1.96 \left( \sqrt{T+1} - \sqrt{T} \right) < \frac{\Delta\bar{y}_{T+1}}{\hat{\sigma}/\sqrt{N}} \quad (18)$$

where the right-hand expression is the marginal expected  $t$ -statistic of the  $T + 1^{th}$  week.

---

<sup>28</sup>Note that this expression is completely general for independent random draws under any marginal indexing or ordering. In the identically distributed case, though, the expected mean for the marginal draw is equal to all inframarginal draws, so the inequality holds.

We can summarize these insights by returning to our formula for the  $t$ -statistic:

$$t_{\Delta\bar{y}_T} = \sqrt{\frac{N}{2}} \left( \frac{\sum_{t=1}^T \Delta\bar{y}_t}{\sqrt{T}\hat{\sigma}} \right).$$

Since the denominator is growing at  $O(\sqrt{T})$ , in order for the  $t$ -statistic to grow, the numerator must grow at a faster rate. In the limit we know this cannot be as the total impact of the advertising would diverge faster than even the harmonic series.<sup>29</sup>

Ex-ante it is hard to know when the trade-off turns against you. The effect may decay slower than the harmonic series initially, and then move towards zero quite quickly. Of course if we knew the pattern of decay, we would have answered the question the whole exercise is asking! So in the end the practitioner must make a judgment call. While choosing longer time frames for advertising effectiveness analyses should capture more of the cumulative effect (assuming that it is generally positive), including additional weeks may just cloud the picture by adding more noise than ad impact. Measuring the effects of advertising inherently involves this sort of “judgment call”—an unsatisfying step in the estimation process for any empirical scientist. But the step is necessary since, as we have shown, estimating the long-run effect of advertising is a losing proposition—the noise eventually overwhelms the signal, the question is “when” and right now our judgment call is to use 1–4 weeks, but this is far from the final word.

---

<sup>29</sup>We note that an asset with infinite (nominal) returns is not implausible per se (a consolidated annuity, known as a “consol,” does this), but we do find infinite effects of advertising implausible. The harmonic series is  $\sum \frac{1}{t}$  whereas the requisite series for an increasing  $t$ -statistic would be  $\approx \sum \frac{1}{\sqrt{t}}$  which diverges much more quickly.

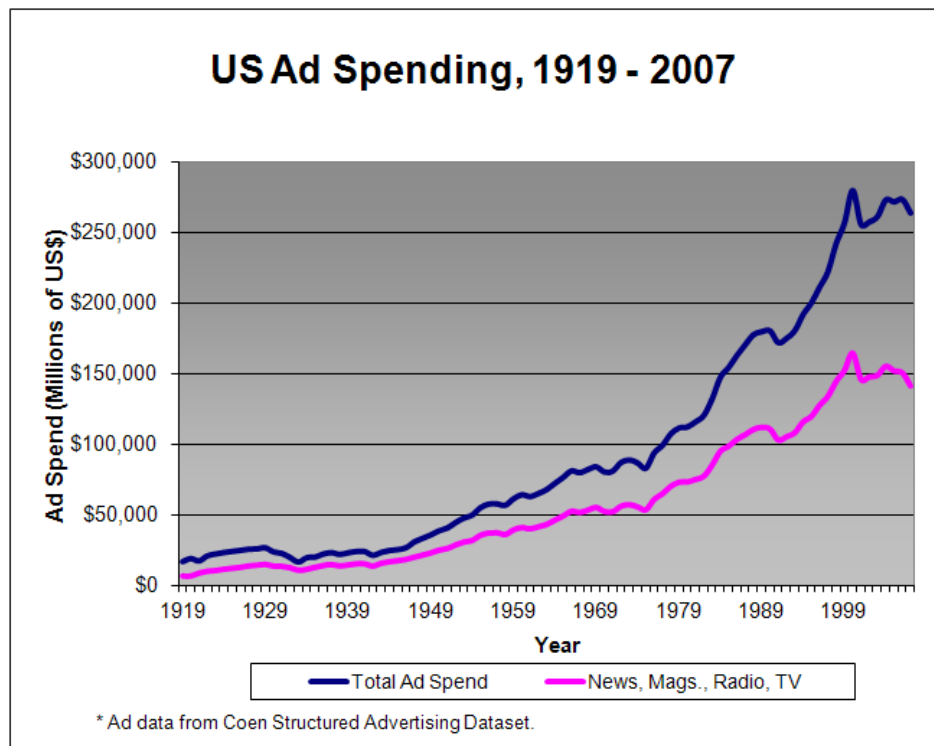


## 6.4 Display ad example



Appendix Figure 1: Display ad example on Yahoo.com.

## 6.5 US ad spending figures



Appendix Figure 2: U.S. Ad Spending 1919–2007.

## 6.6 Advertising across industries and firms

Appendix Table 1: Advertising Expenditure Across Industries and Firms

Industry/Firm	Revenue In \$Billion	Gross margin %	Ad Expenditure In \$Billion	Ad Revenue Share %
<b>Mobile Carriers</b>				
Verizon	114.2	56.9%	1.56344	1.37%
Sprint Nextel	35.1	41.8%	0.67308	1.92%
ATT	127.4	54.5%	1.73602	1.36%
T-Mobile	19.2	N/A	0.52627	2.75%
<b>Automakers</b>				
Honda	115.1	21.4%	0.57124	0.50%
Toyota	262.2	10.2%	0.85032	0.32%
Ford	133.3	17.2%	0.87670	0.66%
GMC	150.1	12.7%	0.17907	0.12%
Fiat-Chrysler	55.0	5.5%	0.87490	1.59%
Hyundai	74.0	N/A	0.30144	0.41%
Dodge	N/A	N/A	0.52501	N/A
<b>Rental Cars</b>				
Avis Budget Group	6.7	24.5%	0.04520	0.67%
Hertz	8.6	43.2%	0.03735	0.43%
Enterprise/Alamo	13.5	N/A	0.06733	0.50%
Dollar Thrifty	1.5	33.7%	0.00021	0.01%
<b>Airlines</b>				
American (AMR)	24.9	47.4%	0.06034	0.24%
United	37.4	56.3%	0.03313	0.09%
Delta	36.5	39.0%	0.05801	0.16%
US Airways	13.7	33.9%	0.01151	0.08%
<b>Online Brokerages</b>				
Scottrade	0.8	100.0%	0.07084	8.45%
Etrade	1.3	100.0%	0.16672	12.63%
TD Ameritrade	2.8	100.0%	0.05034	1.82%
<b>Fast Food</b>				
McDonald's	27.4	39.0%	0.95926	3.50%
Burger King	2.3	37.6%	0.29712	12.92%
Wendy's	2.4	25.3%	0.27248	11.21%
Dairy Queen	2.5	N/A	0.07276	2.91%
Jack in the Box	2.2	45.2%	0.07253	3.30%