

# Measuring the Effects of Advertising: The Digital Frontier\*

Randall Lewis  
Google, Inc.

Justin M. Rao  
Microsoft Research

David H. Reiley  
Google, Inc.

May 29, 2013

## Abstract

Online advertising offers unprecedented opportunities for measurement. A host of new metrics, clicks being the leading example, have become widespread in advertising science. New data and experimentation platforms open the door for firms and researchers to measure true causal effects of advertising on a variety of consumer behaviors, such as purchases. We dissect the new metrics and methods currently used by industry researchers, attacking the question, “How hard is it to reliably measure advertising effectiveness?” We outline the questions that we think can be answered by current data and methods, those that we believe will be in play within five years, and those that we believe could not be answered with arbitrarily large and detailed data. We pay close attention to the advances in computational advertising that are not only increasing the impact of advertising, but also usefully shifting the focus from “who to hit” to “what do I get.”

## 1 Introduction

In the United States, advertising is a \$200 billion industry, annually. We all consume “free” services monetized by consumer attention to advertising including network television, email, social networking, and a vast array of news and other content on the World Wide Web. Yet despite representing a relatively stable 2% of GDP since World War I and paying for activities that comprise most of Americans’ leisure time (American Time Use Survey, 2010), advertising remains poorly understood by economists—historically data have been insufficient to measure the true impact of advertising on consumer behavior. Theories of advertising that have important implications for competition are even harder to empirically validate. The digital era offers an unprecedented opportunity to bridge this informational divide. The potential advances can be attributed to two key factors: 1) ad delivery and purchases at the individual level can be linked and made available to advertisers at low cost, and 2) ad delivery can be randomized, generating exogenous variation essential to identifying causal effects. It is hard to understate the opportunities over traditional media created by the confluence of these two factors.

---

\*Much of this work was done when all the authors were at Yahoo! Research. We thank Garrett Johnson, Dan Nguyen, Sergiy Matusevych, Iwan Sakran, Taylor Schreiner, Valter Sciarillo, Christine Turner, Michael Schwarz, Preston McAfee, and numerous other colleagues for their assistance and support in carrying out the research.

In traditional (offline) settings, the econometric measurements of advertising effectiveness typically rely on aggregate data fraught with endogeneity and identification problems.<sup>1</sup> The ease of data collection in online advertising has led to standard reporting of quantitative data for advertising campaigns, most notably the click-through rate (CTR). Of course, the CTR of an ad is only an intermediate proxy for the real outcome of interest to the advertiser: increased purchases by consumers, both in the present and future (perhaps through increased “brand awareness” today). In some cases advertisers can also obtain a “conversion rate,” the ratio of transactions attributed to the campaign to ad exposures. This measure seems ideal, but the attribution step is critical and current methods of assigning attribution have serious flaws, which we discuss in detail.

Building on this point, we have discovered a number of conceptual flaws in standard industry data collection and analysis methods used to measure the effects of advertising. In other words, the deluge of data on advertising exposures, clicks, and other associated outcomes have not necessarily created greater understanding of the basic causal effects of advertising, much less an understanding of more subtle questions such as the relative effectiveness of different types of consumer targeting, ad creatives, or frequency of exposure. The voluminous data, it seems to us, have created ample opportunity for mistaken inference.

First, many models assume that if you do not click on the ad, then the ad has no effect on your behavior. Here we discuss work by coauthors Lewis and Reiley that showed online ads can drive offline sales, which are typically not measured in conversion or click rates; omitting these non-click-based sales leads to underestimating the total effects of advertising.

Second, many models assume that if you do click on an ad and subsequently purchase, that conversion must have been *due to that ad*. This assumption seems particularly suspect in cases, such as search advertising, where the advertising is deliberately targeted at those consumers most likely to purchase the advertised product and temporally targeted to arrive when a consumer is performing a task related to the advertised good. For example, a person searching for “rental car” is relatively likely to make an online reservation even in the absence of sponsored search advertising. Without a control group to establish the “baseline conversion rate,” these models effectively assume that any purchases must be *caused* by the ad and thus tend to overestimate the effects of advertising.<sup>2</sup>

Third, more sophisticated models that do compare exposed to unexposed users to establish a baseline purchase rate typically rely on natural, endogenous advertising exposure and can easily generate biased estimates due to unobserved heterogeneity (Lewis et al., 2011). This occurs when the pseudo-control group does not capture important characteristics of the treated group, such as purchase intent or browsing intensity, which we show can easily be correlated with purchases

---

<sup>1</sup>A notable exception are the split cable TV experiments reported in Abraham et al. (1995). The sample sizes in these experiments, run in a small U.S. town, were far smaller than online experiments. Bagwell (2005) surveys many studies using the observational methods.

<sup>2</sup>A recent paper from eBay Research Labs addresses this issue with experiments and finds that observational methods severely overstate the impact of some search keywords targeted by the company (Blake et al., 2013).

whether advertising is present or not. Using data from 25 large experiments run at Yahoo! (Lewis and Rao, 2013), we have found that the standard deviation of purchases is typically ten times the mean. With such a noisy dependent variable, even a tiny amount of endogeneity can severely bias estimates. Beyond inducing bias in coefficient estimates, these specification errors also give rise to an over-precision problem. Because advertising typically explains only a very small fraction of the variance in consumer transaction behavior, even cleanly designed experiments typically require over a million subjects in order to be able to measure economically meaningful effects with any statistical precision (but even experiments with 1 million subjects can have surprisingly weak power, depending on the variance in sales).

Since experiments are generally considered the gold standard for precision (treatment is exogenous and independent across individuals), we should be suspicious if observational methods claim to offer higher precision. Further, with non-experimental methods, omitted heterogeneity or selection bias (so long as it can generate a partial R-squared of 0.00005 or greater) can induce bias that swamps plausible estimates of advertising effectiveness. Thus, if an advertiser does not use an experiment to evaluate advertising effectiveness, she has to have a level of confidence in her model that, frankly speaking, we find unreasonable given the obvious selection effects due to ad targeting and synchronization of advertising with product launches (ex. new iPad release) and demand shocks (such as the holiday shopping season).

In this chapter, we report on research involving clean identification through both natural experiments and controlled field experiments. While we feel this work advances the state-of-the art for understanding the economics of online advertising, it has also given us a deeper appreciation for the limits of current data and methods. For example, we show that seemingly simple “cross-channel” complementarity measures (for example, should I run a search ad during a display-advertising campaign?) are exceedingly difficult to reliably estimate. Here we present evidence taken from Lewis and Nguyen (2013) that display advertising can increase keyword searches for the advertised brand. Clicks on sponsored links in this setting could be (incorrectly) attributed entirely to the search ad. But while we can document the effect, the results are statistically noisy, and we cannot tell if search ads perform better or worse when paired with display advertising. In principle one could use our experimental design with much more data to answer this sort of question, but reaching five million individuals may be out of reach<sup>3</sup> for most advertisers. We also note that this type of complementarity question appears to us to be nearly unanswerable for traditional media such as television and billboard, given the practical hurdles in measuring and randomizing delivery on an individual level. Simply put, if you cannot precisely control what people see, you can’t measure what advertising does.

While we have tried to draw attention to questions that have recently become possible to answer with new data sources, we also believe that some questions are still outside the statistical power

---

<sup>3</sup>Pun intended.

of current experimental infrastructure (and thus data) and methods. One example is the long-run effects of advertising. Essentially any analysis of the impact of advertising has to make a judgment call on which time periods to use in the analysis. Often this is the “campaign window” (the time the ads were running) or the campaign window plus a chosen interval of time (we generally used 1-4 weeks). Yet anyone doing this knows it is “wrong” because any impact that occurs after the cutoff should count in the return on investment (ROI) calculation. We explain why practitioners typically choose relatively short impact windows. The intuition is that the longer the time window under study, the lower the signal-to-noise ratio in the data (presuming the ad gets less impactful over time): point estimates of the cumulative effect tend to increase with longer time horizons, but standard errors of the effect increase by even more. This leads to an estimation “impossibility” analogous to the well-known “curse of dimensionality.”

In the final two sections, we discuss how computational methods have increased advertising effectiveness and also helped solve the targeting problem. With automated targeting, the conversation is usefully shifted from “who to hit” to “what should I get.” However, the key parameters of the automated system, such as a bid or valuation of an action such as a click or conversion, the budget of the campaign and the duration, still must be entered by a human. Indeed these are the exact parameters that we have argued are very difficult to estimate. In second to last section, we discuss how new advances in ad-delivery, measurement and infrastructure are creating opportunities to advance the science of advertising. In the final section we present concluding remarks.

## 2 Selection and power

On a daily basis, the average American encounters 25–45 minutes of television commercials (Source: Kantar Media), as well as hundreds of print, radio, billboard, and Internet ads. According to the Coen Structured Advertising Dataset, advertising has accounted for 1.5–2% of GDP since World War I. In today’s dollars this amounts to about \$500 per American per year.<sup>4</sup> So to break even, the universe of advertisers needs to net about \$1.35 in *marginal profits* per person per day. Given the gross margins of firms that advertise, our educated guess is that this roughly corresponds to about \$4-6 in incremental sales.

When an advertiser enters this fray, it must compete for consumers’ attention. The cost per person of a typical campaign is quite low. Online “display” (banners, rectangular units, etc.) campaigns that deliver a few ads per day to a targeted individual cost about 1–2 cents per person per day. Television ads delivered once per person per day are only a bit more expensive. Note that even an aggressive campaign will typically only garner a small percentage of an individual’s daily

---

<sup>4</sup>Mean GDP per American is approximately \$50,000 in 2011, but median household income is also approximately \$50,000. The average household size is approximately 2.5, implying an individual’s share of median household income is roughly \$20,000. Thus, while 2% of GDP actually implies a per capita expenditure of \$1,000, we use \$500 as a round and conservative figure that is more representative of the average American’s ad exposure.

advertising exposure. Stated informally: since ads are cheap, the impact of each ad view should be small. We see many ads per day; only a minority of them are relevant enough to a given person to impact his behavior. The relatively modest average impact per person, in turn, makes it difficult to assess cost-effectiveness.

What complicates matters is that individual-level sales are quite volatile for many advertisers. An extreme example is automobiles—the sales impact is either tens of thousands of dollars, or it is zero.<sup>5</sup> While not as extreme, many other heavily advertised categories, including consumer electronics, clothing and apparel, jewelry, air travel, banking, and financial planning also have volatile consumption patterns.<sup>6</sup> Here we summarize work presented in Lewis and Rao (2013), which used 25 large advertising field experiments to quantify how individual expenditure volatility impacts the power of advertising effectiveness (hereafter, *adfx*) experiments. In general, the signal-to-noise ratio is much lower than we typically encounter in economics.

We now introduce some formal notation to clarify the argument. Consider an outcome variable  $y$  (sales), an indicator variable  $x$  equal to 1 if the person was exposed to the advertising, and a regression estimate  $\hat{\beta}$ , which gives the average difference between the exposed (E) and unexposed (U) groups. In an experiment, exposure is exogenous—determined by a flip of the proverbial coin. In an observational study, one would also condition on covariates  $W$ , 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. We consider a regression of  $y$  on  $x$ , whose coefficient  $\hat{\beta}$  will give us a measure of the average dollar impact of the advertising per consumer.

We use standard notation for 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. We assume for simplicity that the exposed and unexposed samples are the same size ( $N_E = N_U = N$ ) as well as equal variances ( $\sigma_E = \sigma_U = \sigma$ ) to simplify the formulas:

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

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

$$\Delta \bar{y} \equiv \bar{y}_E - \bar{y}_U \quad (3)$$

$$\hat{\sigma}_{\Delta \bar{y}} \equiv \sqrt{\frac{\hat{\sigma}_E^2}{N_E} + \frac{\hat{\sigma}_U^2}{N_U}} = \sqrt{\frac{2}{N}} \cdot \hat{\sigma} \quad (4)$$

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 explained by the advertising (or, in the model with

<sup>5</sup>The marginal profit impact is large, but clearly smaller, as it is the gross margin times the sales impact.

<sup>6</sup>For a bank, the consumption pattern once you sign up might be predictable, but the bank is making money from consumer switching which is “all or nothing.”

covariates, the partial  $R^2$  after first partialling out covariates—for more explanation, see Lovell, 2008):

$$R^2 = \frac{\sum_{i \in U} (\bar{y}_U - \bar{y})^2 + \sum_{i \in E} (\bar{y}_E - \bar{y})^2}{\sum_i (y_i - \bar{y})^2} = \frac{2N \left(\frac{1}{2}\Delta\bar{y}\right)^2}{2N\hat{\sigma}^2} = \frac{1}{4} \left(\frac{\Delta\bar{y}}{\hat{\sigma}}\right)^2 \quad (5)$$

Second is the  $t$ -statistic for testing the hypothesis that the advertising had no impact:

$$t_{\Delta\bar{y}} = \frac{\Delta\bar{y}}{\hat{\sigma}_{\Delta\bar{y}}} = \sqrt{\frac{N}{2}} \left(\frac{\Delta\bar{y}}{\hat{\sigma}}\right) \quad (6)$$

In both cases, we have related a standard regression statistic to the ratio between the average impact on sales and the standard deviation of sales between consumers.

In the following hypothetical example, we calibrate values using approximately median values from 19 retail sales experiments run at Yahoo!. For expositional ease, we will discuss it as if it is a single experiment. The campaign goal is a 5% increase in sales during the two weeks of the campaign, which we will use as our “impact period” of interest. During this period, customers of this advertiser make purchases with a mean of \$7 and a standard deviation of \$75.<sup>7</sup> The campaign costs \$0.14 per customer, which amounts to delivering 20–100 display ads at a price of \$1–\$5 CPM,<sup>8</sup> and the gross margin (markup over cost of goods sold, as a fraction of price) is assumed to be about 50%.<sup>9</sup> A 5% increase in sales equals \$0.35 per person, netting profits of \$0.175 per person. Hence, the goal for this campaign is to deliver a 25% return on investment (ROI):  $\$0.175/\$0.14 = 1.25$ .<sup>10</sup>

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

$$R^2 = \frac{1}{4} \cdot \left(\frac{\$0.35}{\$75}\right)^2 = 0.0000054 \quad (7)$$

That is, even for a *successful* campaign with a *relatively large* ROI, we expect an  $R^2$  of only *0.0000054*. This will require a very large  $N$  to identify any influence at all of the advertising, 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 a ratio of 0.0047 between impact size and standard deviation of sales, the expected  $t$ -stat 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, as the normally distributed  $t$ -stat should be less than the critical value of 1.65 about 5% of

---

<sup>7</sup>Based on data-sharing arrangements between Yahoo! and a number of advertisers spanning the range from discount to high-end retailers, the standard deviation of sales is typically about 10 times the mean. Customers purchase goods relatively infrequently, but when they do, the purchases tend to be quite large relative to the mean.

<sup>8</sup>CPM is the standard for impression-based pricing for online display advertising. It stands for “cost per mille” or “cost per thousand;” M is the roman numeral for 1,000.

<sup>9</sup>We base this assumption on our conversations with retailers and our knowledge of the industry.

<sup>10</sup>For calibration purposes, note that if the gross margin were 40% instead of 50%, this would imply a 0% ROI.

the time given the true effect is a 25% ROI. With 200,000 unique customers, the expected t-stat is 1.04, indicating the test is hopelessly underpowered to reliably detect an economically relevant impact: under the alternative hypothesis of a healthy 25% ROI, we fail to reject the null 74% of the time.<sup>11</sup>

The low  $R^2 = 0.0000054$  for the treatment variable  $x$  in our hypothetical randomized trial has serious implications for observational studies, such as regression with controls, difference-in-differences, and propensity score matching. A very small amount of endogeneity would *severely bias* estimates of advertising effectiveness. 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. Compare this to a classic economic example such as the Mincer wage/schooling regression (Mincer, 1962), in which the endogeneity is roughly 1/8 the treatment effect (Card, 1999). For observational studies, it is always important to ask, “What is the partial  $R^2$  of the treatment variable?” If it is very small, as in the case of advertising effectiveness, clean identification becomes paramount, as a small amount of bias can easily translate into an economically large impact on the coefficient estimates.

Our view has not yet been widely adopted, however, as evidenced by the following quotation from the president of comScore, a large data-provider for online advertising:

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 argument we have made shows that simply comparing exposed to unexposed can lead to bias that is many orders of magnitude larger than the true size of the effect. Indeed, this methodology led the author to report as much as a 300% improvement in outcomes for the exposed group, which seems surprisingly high (it would imply, for instance, that advertisers are grossly underadvertising). Since all ads have some form of targeting,<sup>12</sup> endogeneity is always a concern. For example, most display advertising aims to reach people likely to be interested in the advertised product, where such interest is inferred using demographics or past online behavior of that consumer. Similarly, search advertising targets consumers who express interest in a good at a particular point in time, where the interest is inferred from their search query (and potentially past browsing behavior). In these cases, comparing exposed to unexposed is precisely the *wrong* thing to do. By creating

---

<sup>11</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.

<sup>12</sup>“Untargeted” advertising usually has implicit audience targeting based on where the ads are shown or implicit complementary targeting due to other advertisers purchasing targeted inventory and leaving the remnant inventory to be claimed by advertisers purchasing “untargeted” advertising inventory.

exogenous exposure, the first generation of advertising experiments have been a step in the right direction. Experiments are ideal—necessary, in fact—for solid identification.

Unfortunately, for many advertised products the volatility of sales means that even experiments with millions of unique users can still be underpowered to answer basic questions such as “Can we reject the null hypothesis that the campaign had zero influence on consumer behavior?” Measuring sales impact, even in the short-run, turns out to be much more difficult than one might have thought. The ability to randomize ad delivery on an individual level and link it to data on customer-level purchasing behavior has opened up new doors in measuring advertising effectiveness, but the task is still by no means easy. In the remainder of the paper we discuss these challenges. The next section focuses on using the right metrics to evaluate advertising.

### 3 The click and related metrics

The click-through-rate, or CTR, has become ubiquitous in the analysis and decision-making surrounding online advertising. It is easy to understand why: clicks are cleanly defined, easily measurable, and occur relatively frequently. An obvious but intuitively appealing characteristic is that an ad-click cannot occur in the absence of an ad. If one runs 100,000 ads and gets a 0.2% CTR (a decent rate for a display ad or a low-ranked search ad), it is tempting to conclude the ad caused 200 new website visits. The assumption may well be approximately true for little-known brands. But for well-known brands, there are important ways that consumers might navigate to the site in the absence of an ad, such as navigating to a favorite online store or finding it in organic (that is, not paid or “sponsored”) search results on a topic like “car rental.”<sup>13</sup> It is a mistake to assume that *all* of those 200 visits would not have occurred in the absence of the ad—that is, those clicks may be crowding out visits that would have happened via other means—but without a control group one has no way to pin down the counterfactual.

While the overcounting problem is important, it can probably be pinned down with randomized trials where the control group is used to estimate the “baseline arrival rate.” For example, a sponsored search ad could be turned off during random times of the day and the firm could measure arrivals from the search engine for when the ad is running and when it is not (this approach is used in Blake et al., 2013).<sup>14</sup> A deeper problem with the CTR is what it misses. First, it does little for “brand advertisers”—firms that are not trying to generate immediate online sales, but rather to promote awareness and good-will for the brand. To assess their spend, brand advertisers have traditionally relied on surveys that attempt to measure whether a campaign raised the opinion of the firm in the minds of their target consumers. Linking the surveys to future purchasing

---

<sup>13</sup>Research at Google has confirmed that ads can crowd out organic clicks ((Kumar and Yildiz, 2011; Chan et al., 2010)).

<sup>14</sup>Despite the simplicity of their design, Blake et al. estimate that their employer, eBay, had been wasting tens of millions of dollars a year.

behavior adds another layer of complexity, both because the time frame from exposure to sale is longer (something we will discuss in more detail in Section 5) and because it requires a reliable link from hypothetical responses to actual behavior, which can be fraught with what is known as “hypothetical bias” (Dickie et al., 1987; Murphy et al., 2005).

Another class of advertisers sell goods both online and in brick-and-mortar stores. Lewis and Reiley (2013b) show that for a major retailer, the *majority* of the sales impact comes offline. Johnson, Lewis, and Reiley (2013) link the offline impact to consumers who lived in close physical proximity to one of the retailer’s locations. The click fundamentally misses any offline impact of an ad, and the evidence indicates that this can introduce a large, negative measurement bias for advertisers that also do business through brick-and-mortar stores or by phone.

The click is an *intermediate* metric for short-term sales (and short-term sales might be considered a further upstream metric for the net-present-discounted value of a customer). Alternatively, however, advertisers can now run “cost per acquisition” (CPA) advertising on many ad exchanges. An acquisition, or conversion, is defined as a successful transaction that has a qualifying connection to the advertisement. On the surface, focusing on conversions seems more attractive than clicks because it is a step closer to sales. Unfortunately this benefit brings with it what is known as the “attribution problem:” which ad gets “credit” for a given sale? Suppose a consumer views and clicks a given ad, but does not purchase on the same day. Over the next few days, she sees a host of other ads for the product (which is likely, given a practice known as “re-targeting”) and then purchases the good. Which ad should get credit for the purchase?

Ad exchanges tend to use a set of rules to solve these problems from an accounting perspective. Common rules include requiring a click for credit or only counting the “last click” (so if a consumer clicks a re-targeted ad, that ad gets credit). Requiring a click seems to make sense and is enormously practical as it means a record of all viewers that see the ad but do not click need not be saved.<sup>15</sup> However, requiring a click errs in assuming that ads can only have an impact through clicks, which is empirically not true (Lewis et al., 2012). The “last click” rule also has intuitive appeal. The reasoning goes as follows: had the last click not occurred, the sale would not have happened. Even if this were true, which we doubt, the first click or ad view might have led to web search, or other activity, including the behavioral markers used for re-targeting, which made the last click possible. The causal attribution problem is typically “solved” by ad-hoc assumptions, such as “the first ad and the last ad viewed before purchase each get 40% of the credit, while the intermediate ad views share the remaining 20% of the credit for the purchase.”<sup>16</sup> A proliferation of such rules gives practitioners lots of choices, but none of them necessarily gives an unbiased measurement of the performance of their ad spending. In the end, such complicated payment rules might make the click more attractive after all.

---

<sup>15</sup>CTRs are commonly  $\approx 0.1\%$ , meaning, storage and processing costs of only clicks involves only  $\frac{1}{1,000}$  of the exposure and click data.

<sup>16</sup>Source: <https://support.google.com/analytics/bin/answer.py?hl=en&answer=1665189>

The attribution problem is also present in the question of complementarities between display and search advertising. Recent work has shown that display ads causally influence search behavior (Lewis and Nguyen, 2013). The authors demonstrate this by comparing the search behavior of users exposed to the campaign ad to users who would have been served the campaign ad but were randomly served a placebo. Brand-related keywords were significantly more prevalent in the treatment group as compared to the control. The attribution problem has received more attention in online advertising because of the popularity of cost-per-acquisition and cost-per-click payment mechanisms, but it applies to offline settings as well. How do we know, for example, whether an online ad was more responsible for an online conversion than was the television ad that same user saw? Nearly every online campaign occurs contemporaneously with a firm’s offline advertising through media such as billboards and television because large advertisers are continuously advertising across many media.<sup>17</sup> Directly modeling the full matrix of first-order interactions is well beyond the current state of the art. Indeed in every paper we know of evaluating online advertising, the interactions with offline spending is ignored.

This criticism applies to our own work as well, but at least with controlled experiments we know we are measuring the true marginal effect of an online campaign while holding all other advertising constant. However, suppose contemporaneous television advertising increases the marginal effectiveness of online display advertising. In that case, the treatment-control comparison captures the effectiveness in the presence of the television spending, which would be an overestimate if one were interested in the effectiveness of the online advertising in the absence of television spending. However, we find the opposite case more plausible: global diminishing returns to advertising would make online advertising less effective in the presence of television ads and would have to be true (at least locally) in order for firms to be near the optimum level of ad spending. If increasing returns or economies of scope in advertising were present at current spending levels, that would mean that firms were currently under-advertising by a large margin. While we do not believe that firms can be exactly at the optimum given the difficulty of measuring effectiveness, we do think that advertising levels sufficiently high to produce diminishing returns and substitutabilities are more likely than not. In this case, online experiments *underestimate* how effective online advertising would be if firms cut their television budgets and switched that spending online. The same of course could be said for television experiments, if they were feasible. Online experiments right now give us one point to work with, so to speak, but global diminishing returns and cross-channel effects make it difficult to know exactly what do with this estimate.<sup>18</sup>

---

<sup>17</sup>(Lewis and Reiley, 2013a) show that Super Bowl commercials cause viewers to search for brand-related content across a wide spectrum of advertisers.

<sup>18</sup>Substituting across ad channels has been shown to impact ad prices (Goldfarb and Tucker, 2011).

## 4 A case study of a large-scale advertising experiment

To get a better idea of how large advertising experiments are actually run, in this section we present a case study taken from Lewis and Reiley (2013b) (herein “LR”). LR ran a large-scale experiment for a major North American retailer. The advance the paper makes is linking existing customers in the retailer’s sales records, for both online and brick-and-mortar sales, to a unique online user identifier, in this case the customer’s Yahoo! username.

The experiment was conducted as follows. The match yielded a sample of 1,577,256 individuals who matched on name and either email or postal address. The campaign was targeted only to existing customers of the retailers as determined by the match. Of these matched users, LR assigned 81% to a treatment group who subsequently viewed two advertising campaigns promoting the retailer when logged into Yahoo’s services. The remaining 19% were assigned to the control group and prevented from seeing any of the retailer’s ads from this campaign on the Yahoo! network of sites. The simple randomization was designed to make the treatment-control assignment independent of all other relevant variables.

Table 1: Summary Statistics for the Campaigns

	Campaign 1	Campaign 2	Both Campaigns
Time Period Covered	Early Fall '07	Late Fall '07	
Length of Campaign	14 days	10 days	
Number of Ads Displayed	32,272,816	9,664,332	41,937,148
Number of Users Shown Ads	814,052	721,378	867,839
% Treatment Group Viewing Ads	63.7%	56.5%	67.9%
Mean Ad Views per Viewer	39.6	13.4	48.3

*Source:* Lewis and Reiley (2013b).

The treatment group of 1.3 million Yahoo! users was exposed to two different advertising campaigns over the course of two months in fall 2007, separated by approximately one month. Table 1 gives summary statistics for the campaigns, which delivered 32 million and 10 million impressions, respectively. The two campaigns exposed ads to a total of 868,000 users in the 1.3-million-person treatment group. These individuals viewed an average of 48 ad impressions per person.

The experiment indicated an increase in sales of nearly 5% relative to the control group during the campaign, a point estimate which would translate to an extremely profitable campaign (with the retailer receiving nearly a 100% rate of return on the advertising spending). However, purchases had sufficiently high variance (due in part to 95% of consumers making zero purchases in a given week) to render the point estimate not statistically significantly different from zero at the 5% level. Controlling for available covariates (age, gender, state of residence) did not meaningfully reduce standard errors. This is good example of how economically important effects of advertising can

be statistically very difficult to detect, even with a million-person sample size. Just as we saw in Section 2, we see here that the effects of advertising are so diffuse, explaining such a small fraction of the overall variance in sales, that the statistical power can be quite low. For this experiment, power calculations show that assuming the alternative hypothesis that the ad broke even is true, the probability of rejecting the null hypothesis of zero effect of advertising is only 21%.

The second important result of this initial study was a demonstration of the biases inherent in using cross-sectional econometric techniques when there is endogenous advertising exposure. This is important because these techniques are often employed by quantitative marketing experts in industry. Abraham (2008), for example, advocates comparing the purchases of exposed users to unexposed users, despite the fact that this exposure is endogenously determined by user characteristics and browsing behavior, which might easily be correlated with shopping behavior. To expose the biases in these methods, LR temporarily “discarded” their control group and compared the levels of purchases between exposed and (endogenously) unexposed parts of the treatment group. The estimated effects of advertising were three times as large as in the experiment, and with the opposite sign! This erroneous result would also have been deemed highly statistically significant. The consumers who browsed Yahoo! more intensely during this time period (and hence were more likely to see ads) tended to buy less, on average, at the retailer, regardless of whether they saw the ads or not (this makes sense, because as we will see most of the ad effect occurred offline). The control group’s baseline purchases prior to the ad campaign showed the same pattern. Without an experiment an analyst would have had no way of realizing the extent of the endogeneity bias (in this case, four times as large as the true causal effect size) and may have come to a strikingly wrong conclusion.

Observing the consistent differences between exposed and unexposed groups over time motivated LR to employ a difference-in-differences estimator. Assuming that any unobserved heterogeneity was constant over time allowed LR to take advantage of both exogenous and endogenous sources of variation in advertising exposure, which turned out to reduce standard errors to the point where the effects were statistically significant at the 5% level. The point estimate was approximately the same as (though slightly higher than) the straight experimental estimate, providing a nice specification check. With this estimator, LR also demonstrated that the effects of the advertising were persistent for weeks after the end of the campaign, that the effects were significant for in-store as well as online sales (with 93% of the effect occurring offline), and that the effects were significant even for those consumers who merely viewed but never clicked the online ads (with an estimated 78% of the effect coming from non-clicking viewers). In a companion paper (Lewis and Reiley, 2012), the authors also showed that the effects were particularly strong for the older consumers in the sample—sufficiently strong to be statistically significant even with the simple (less efficient) experimental estimator.

In a follow-up study, Johnson, Lewis, and Reiley (2013, henceforth JLR) improved on some of the

weaknesses of the design of the original LR experiment. First, JLR ran “control ads” (advertising one of Yahoo’s own services) to the control group, allowing them to record which control-group members would have been exposed to the ad campaign if they had been in the treatment group. This allowed them to exclude from their analysis those users (in both treatment and control groups) who were not exposed to the ads and therefore contributed noise but no signal to the statistics. Second, JLR convinced the advertiser to run equal-sized treatment and control groups, which improved statistical power relative to the LR article’s 81:19 split. Third, JLR obtained more detailed data on purchases: two years of pre-campaign sales data on each individual helped to explain some of the variance in purchases, and disaggregated daily data during the campaign allowed them to exclude any purchases that took place before the first ad delivery to a given customer (which therefore could not have been caused by the ads, so including those purchases merely contributed noise to the estimates). The more precise estimates in this study corroborate the results of LR, showing point estimates of a profitable 5% increase in advertising which are statistically significant at the 5% level, though the confidence intervals remain quite wide.

## 5 Activity bias

We believe it is not possible to reliably estimate the causal effects of advertising using observational methods. In the preceding sections, we have presented this argument on an abstract level, arguing that the since the partial  $R^2$  of advertising, even for a successful campaign, is so low (on the order of 0.00001 or less), the likelihood of omitted factors not accounting for this much variation is unlikely, especially since ads are targeted across time and people. In this section we show that our argument is not just theoretical. Here identify a bias that we believe is present in most online ad serving; in past work we gave it the name “activity bias” (Lewis et al., 2011). “Activity bias” is a form of selection bias based on the following two features of online consumer behavior: 1) since one has to be browsing online to see ads, those browsing more actively on a given day are more likely to see your ad, and 2) active browsers tend to do more of *everything online*, including buying goods, clicking links and signing up for services. In a non-experimental study, the unexposed group, as compared to the group exposed to an ad, typically failed to see the ad for one or both of the following reasons: the unexposed users browsed less actively or the user did not qualify for the targeting of the campaign. When the former fails, we have activity bias. When the latter fails, we have classic selection bias.

In our 2011 paper, we explored three empirical examples demonstrating the importance of activity bias in different types of web browsing. The first application investigates the causal effects of display ads on users’ search queries. In Figure 5 we plot the time series of the number of searches by exposed users for a set of keywords deemed to be brand-relevant for a firm. The figure shows results for a time period that includes a one-day display-advertising campaign for a national brand on [www.yahoo.com](http://www.yahoo.com).

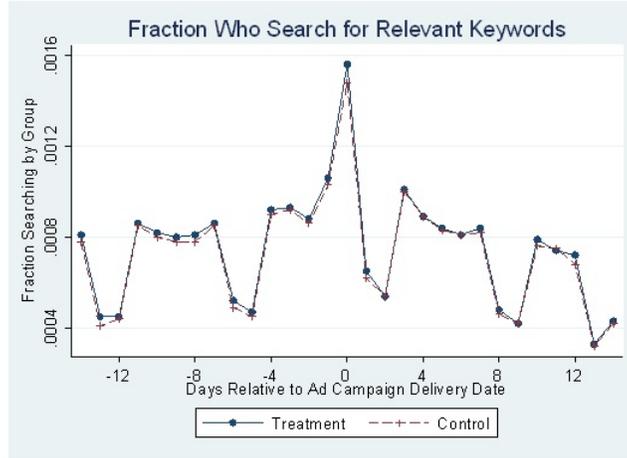


Figure 1: Brand keyword search patterns over time. *Source: Lewis, Rao and Reiley (2011)*

The campaign excluded a randomized experimental control group, though for the moment we ignore the control group and focus on the sort of observational data typically available to advertisers (the treatment group, those that saw the firm’s advertisements). The x-axis displays days relative to the campaign date, which is labeled as Day 0. One can easily see that on the date of the ad, ad viewers were much more likely to conduct a brand-relevant search than on days prior or following. The advertising appears to *double* baseline search volume. Is this evidence of a wildly successful ad? Actually, no. Examining the control group, we see almost the same trend. Brand-relevant keyword searches spike for even those who saw a totally irrelevant ad. What is going on? The control group is, by design of the experiment just as active online as the treatment group, searching for more of *everything*, not just the brand-relevant keywords of interest. The time series also shows that search volume is positively serially correlated over time and shows striking day of week effects—both could hinder observational methods. The true treatment-control difference is a statistically significant, but far more modest, 5.1%. Without an experiment, we would have no way of knowing the baseline “activity-related increase” which we infer from the control group. Indeed, we might have been tempted to conclude the ad was wildly successful.

Our second application involves correlation of activity not just across a publisher and search engine, but across very different domains. We ran a marketing study to evaluate the effectiveness of a video advertisement promoting the Yahoo! network of sites. We recruited subjects on Amazon Mechanical Turk, showed them the video and gave them a Yahoo! cookie so we could track their future behavior. Using the cookie we could see if the ad really generated more Yahoo! activity. The control group saw a political ad totally unrelated to Yahoo! products and services. Again, we ignore the control group to begin. Figure 5 has the same format as Figure 5 Day 0 on the x-axis labels the day an individual saw the video ad (with the actual calendar date depending on the day the subject participated in the study).

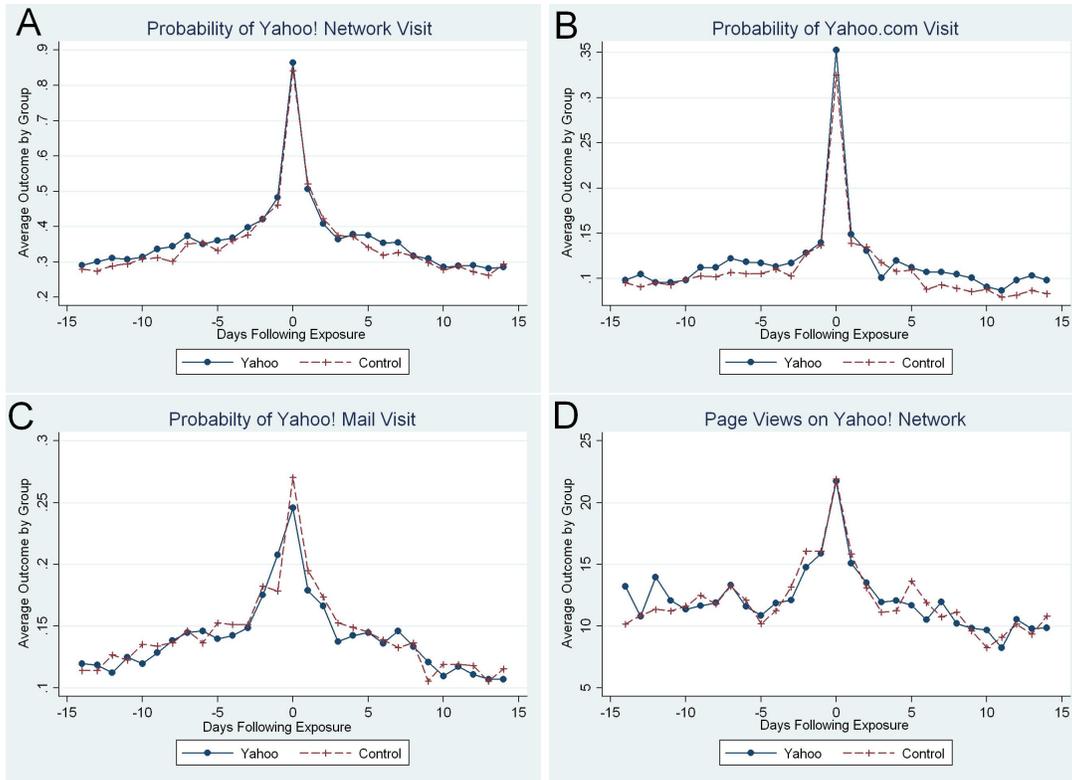


Figure 2: The effect on various Yahoo! usage metric of exposure to treatment/control ads. Panels A-C: Probability of at least 1 visit to the Yahoo! network, Yahoo.com and Mail respectively. Panel D: Total page views on the Yahoo! network. *Source: Lewis, Rao and Reiley (2011)*

Examining the treatment group, we can see that on the day of and the days following ad-exposure, subjects were much more likely to visit a Yahoo! site as compared to their baseline propensity, indicating a large apparent lift in engagement. However, data on the control group reveals the magnitude of activity bias—a very similar spike in activity on Yahoo! occurs on the day of placebo exposure as well. Both groups also show some evidence of positive serial correlation in browsing across days: being active today makes it more likely that you will be active tomorrow as compared to several days from now. People evidently do not engage in the same online activities (such as visiting Yahoo! and visiting Amazon Mechanical Turk) every day, but they engage in somewhat bursty activity that is contemporaneously correlated across sites. Online activity leads to ad exposure, which mechanically tends to occur on the same days as outcome measures we hope to affect with advertising. In the absence of a control group, we can easily make errors in causal inference due to activity bias. In this particular case, the true causal effect of the ad was estimated to be small and not statistically significant—given the cost of running a video ad, it was probably not worth showing, but the biased estimates would have led us to a wrong conclusion in this regard.

The third application again involves multiple websites. This time the outcome measure was

filling out a new-account sign-up form at an online brokerage advertised on Yahoo! Finance. Again our results show that even those who were randomly selected to see irrelevant placebo ads were much more likely to sign up on the day they saw the (placebo) ad than on some other day. We refer the reader to our original paper for the details, stating here that the results are very similar to the ones we have just presented (the now familiar mountain-shaped graphs are again present). With activity bias it seems that one could erroneously “show” that nearly any browsing behavior is caused by nearly any other browsing behavior! We hope that our results will cause industry researchers to be more cautious in their conclusions. Activity bias is a real form of bias that limits the reliability of observational methods.

In the absence of an experiment, researchers may be able to use some other cross-validation technique in order to check the robustness of causal effects. For example, one could measure the “effect” of movie advertisements on searches for the seemingly irrelevant query “car rental.” Similarly, one could check whether (placebo) ad views of a Toyota ad on the New York Times website on May 29 “causes” the same effect on Netflix subscriptions that day as did the actual Netflix ad on the New York Times website on May 30. Differences in differences using such pseudo-control groups will likely give better estimates of true causal effects than simple time-series or cross-sectional studies, though of course a randomized experiment is superior if it is available (Lewis et al., 2011).<sup>19</sup>

Is activity bias a new phenomenon that is unique to the online domain? While it is not obvious that offline behavior is as bursty and as contemporaneously correlated as online behavior, before our study we did not think these patterns were obvious in online behavior either (and scanning industry white papers, one will see that many others still do not find it obvious!). We believe the importance of activity bias in the offline domain is an open question. It is not difficult to come up with examples in which offline advertising exposure could spuriously correlate with dependent variables of interest. Billboards undoubtedly “cause” car accidents. Ads near hospitals “cause” illness. Restaurant ads near malls probably “cause” food consumption in general. Exposure to ads in the supermarket “saver” are likely correlated with consumption of unadvertised products. And so forth. The superior quality of data (and experiments) available in online advertising has laid bare the presence of activity bias in this domain. We believe the level of activity bias in other domains is an interesting, open question.

## 6 Long-run effects

Any study of advertising effectiveness invariably has to specify the window of time to be included in the study. While effects of advertising could in principle last a long time, in practice we must pick

---

<sup>19</sup>In some cases, even such placebo tests may fail as the qualifications for seeing the ad may be intrinsically correlated with the desired outcome as may be the case for remarketing and other forms of targeting which account for search activity and browsing behavior.

a cut-off date, unless marketing officers wish to suffer the fate of Keynes’ long-run economists.<sup>20</sup> From a business perspective, making decisions quickly is an asset worth trading (some) decision accuracy for. But can patient scholars (or firms) hope to measure the long-run effects of advertising? Here we address the statistical challenges of this question. The answer, unfortunately, is rather negative. As one moves further and further from the campaign date, the cumulative magnitude of the sales impact tends to increase. (This is not guaranteed, as ads could simply shift purchases forward in time, so a short time window could measure a positive effect while a long time window gives a zero effect. But in practice, we have so far noticed point estimates of cumulative effects to be increasing in the time window we have studied.) However, the amount of noise in the estimate tends to increase faster than the increase in the signal (treatment effect) itself because in the additional data the control and treatment groups look increasingly similar, making long-run studies less statistically feasible than short-run ones. In the remainder of this section we formalize and calibrate this argument.

We again employ the  $t$ -statistic from above, but include indexes 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>21</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:

$$t_{\Delta\bar{y}_T} = \sqrt{\frac{N}{2}} \left( \frac{\sum_{t=1}^T \Delta\bar{y}_t}{\sqrt{T}\hat{\sigma}} \right). \quad (8)$$

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

---

<sup>20</sup>John Maynard Keynes is quoted as saying, “In the long run, we’re all dead.”

<sup>21</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.

condition necessary for an additional week to increase the  $t$ -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

$$\frac{1}{2} \cdot \frac{1}{T} \sum_{t=1}^T \Delta\bar{y}_t < \Delta\bar{y}_{T+1}. \quad (9)$$

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 easily cloud the previous weeks, as its signal-to-noise ratio is not sufficiently large enough to warrant its inclusion.<sup>22</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

$$1.96 \left( \sqrt{T+1} - \sqrt{T} \right) < \frac{\Delta\bar{y}_{T+1}}{\hat{\sigma}/\sqrt{N}} \quad (10)$$

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

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  $O\left(\sqrt{T}\right)$ , 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>23</sup>

<sup>22</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 always holds.

<sup>23</sup>We note that an asset with infinite (nominal) returns is not implausible per se (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

Now 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.

## 7 Local vs. Global Optimization

Throughout this paper we have discussed the importance of targeting across people (based on preferences) and time (based on “session intent” or demand seasonality). In traditional media, targeting is typically a human-controlled process of determining the demographic groups most likely to consume the product. Readers may be familiar with Nielsen ratings for television, which break down to viewership by demographic categories. These categories are then used to target viewers based on past purchasing data, focus groups and so on. For instance, a show that draws mainly 18–35 year-old males might be well-suited for a shaving commercial. Campaigns often have “reach goals” for specific demographics; marketing representatives use a portfolio of media outlets to meet these goals.

Online advertising is different. Because online systems both gather information about specific users and make the ad-serving decision in *real-time*, the field of “computational advertising” was born. According to one of the founders of the field, Andrei Broder, computational advertising is “a principled way to find the best match between a given user in a given context and a suitable advertisement” (Broder, 2008). In traditional media, you have to specify who you want to advertise to. With computational advertising, you instead specify what you want, an end-goal supported by the system, and automated systems determine how to achieve that goal most efficiently. The end goal could be online sign-ups, clicks to a sales page, and so on. A system will not support all end goals, and some supported goals, such as conversions, might exhibit slow learning because the success rate is so low (1 in 300,000 would not be uncommon for account sign-ups, for instance).

While the details of these systems are well beyond the scope of this paper, we shall give a small taste here. Which display ad to show can be modeled as a multi-armed bandit problem. The possible ads are the “arms” and a user-ad pair is a “pull of the arm.” Papers in this literature

---

$t$ -statistic would be  $\approx \sum \frac{1}{\sqrt{t}}$  which diverges much more quickly.

adapt classic machine learning tools to the ad-serving context (see for instance Pandey et al., 2007). Another approach, which is used in text-link based advertising, is to view the advertisement as a document that must be retrieved and matched to the page the ad is served on, which can be thought of as the query. As indicated by the terminology, this approach models ad-serving using the tools of web search (Rusmevichientong and Williamson, 2006; Cary et al., 2007).

We view (the current incarnation of) computational advertising essentially as automated targeting. It helps *locally* optimize ad spend by minimizing costs for given a campaign goal. Using a high-level goal, such as clicks or sign-ups, combined with a budget, it reduces the need to set targeting dimensions (reduces not eliminates because one might still set priors for the learning system, which might matter a lot in slow-to-learn tasks). Focusing on high-level goals also helps shift the conversation from “who should get ads” to “what do we want to get from our ad spending.” Practitioners should be cautious, however, that the system does not conflate “the audience most likely to convert” to the “audience that delivers the most *additional* conversions.” To see the difference, imagine a customer that would buy anyway, but finds it convenient to click on an ad if he sees one. Paying for this conversion is a total waste of money. In our experience, some automated systems fail to draw this distinction and in doing so, “order anticipate” by advertising to people likely to make a future purchase anyway.

Overall we think computational advertising is promising advance. Conceptually, focusing on end goals can help practitioners manage spend. But it is important not to overstate what local optimization gives you. It does deliver cheaper/quicker clicks or conversions when the machine targeting systems can beat humans, but it does not say how much you should be advertising, how much you should bid per action, or how much of an given action should be attributed to a given ad—the very parameters that determine ad spending and profitability. For instance, suppose an online brokerage calculates that it nets \$100 in profit from every account sign-up. Should it specify \$100 as a maximum bid an automated system and then “set it and forget it?” Well presumably the brokerage is advertising heavily on TV and other media, including other online media that was not the “last click.” Bidding \$100 effectively says all this other spending gets zero credit — the firm would over-advertise using this rule. Of course this is just the “attribution problem” reframed from the advertiser’s perspective.

The point can be stated succinctly: most of the difficulties we have discussed about globally optimizing ad spending apply to computational advertising as well. The one difficulty that does not is determining who to target (locally optimizing). This is a big step forward for practitioners and moving the discussion from “who” to “what” has changed the culture of advertising, but it is not a magic bullet. One could imagine a system that conducts automated experiments to measure incremental conversions and self-governs bids based on the experimental feedback, but we are not there yet in practice.

## 8 Moving forward

Thus far we have highlighted how digital data has opened up many doors in measuring advertising effectiveness, but we have also drawn attention to the challenges still present. In this section we look toward the future and discuss how we think many of the existing challenges will be overcome. Overall we expect the advances to mainly come from better experimentation infrastructure and thus better data, and not necessarily better data analysis (the beauty of randomized trials is that the analysis is simple).

The first advance we will discuss involves reducing the cost of experimentation. The first generation of field experiments we ran at Yahoo! randomly selected a relatively small sample of users targeted by the campaign to see an unrelated advertisement. The problem was that an unrelated ad had to be entered into the booking system and run for the users that were randomized into the control group. The booking system was set up so that a firm could run multiple “creatives” (different versions of the ad) for fractions of traffic. What this meant was the unrelated ad had to be fully specified in the booking system. The firm for whom we ran the experiment did not want to let another retailer get the traffic, because the competitor would benefit from the targeting dimensions set up by the retailer (including, for instance, past purchasing behavior). The solution was to use charity ads for the control group. But this meant that either the advertiser had to pay for the control ads or Yahoo! had to donate them—both options came at a cost that increased linearly in the size of the control group, meaning that first generation experiments had relatively small control groups.

A small control group not only hurts power but also makes experimentation less useful as an evaluative tool. An experiment with 90% of subjects in the treatment group and 10% in the control has the same power as one with 10% in the treatment and 90% control. If control ads are free, then the an advertiser could run 9 of the latter for the cost of 1 of the former.<sup>24</sup> For control ads to be free, the ad server needs to be able to serve the “next ad in line” every time a user is randomized into the control group. Technologically this requires a short serving latency between the request to the ad server, the randomization, and the request for the replacement ad. The replacement ads are known as “ghost ads”—ads that naturally qualified to be served to a given user targeted by the campaign under study but not associated with the advertiser. Ghost ads make exploration and evaluation cheaper. Small treatment groups limit cost and allow advertisers to hone copy early in a campaign, while free control subjects help evaluate the campaign ex-post.

Major online publishers are developing similar experimentation platforms. As experiments become cheaper and easier to run, advertisers will be able to form more precise beliefs on effectiveness

---

<sup>24</sup>Note that the statistical gains from such a change in experimental design are 3-fold. Further altering the design, assuming constant returns to scale from advertising (Lewis, 2010; Johnson et al., 2012), by concentrating the 90% treatment group’s ad impressions all within a smaller 10% treatment group expects an impact that is 9 times as large, resulting in the equivalent ad effectiveness insights from running 81 of the 90%/10% experiments, producing confidence intervals of the ROI that are 9 times more precise at no additional advertising cost.

than has heretofore been possible. Indeed campaigns could be evaluated using an informative prior, which would help combat the power concerns we detailed earlier. With repeated experiments, beliefs would converge to the true effectiveness of the campaign.

Another experimentation technology that improves power is the pre-experiment matching of users. To see how this works, consider an experiment with subjects spread across treatment and control 50-50. A standard experiment would simply flip a coin each time a user arrived at the website and show the ad corresponding to the outcome of the flip. Matching works as follows. Specify a set of attributes you care about such as recent sales. Form pairs of users by minimizing some objective function that defines the distance between two nodes in the graph of users. Then for each pair, flip a coin to determine experimental grouping. By construction the specified metrics should be almost exactly equal between the two groups. For evaluating a noisy variable such as sales, guaranteeing the pre-period sales were the same can be useful. The treatment assignment is still totally exogenous, so all our normal intuition on how experiments identify causal effects goes through. In large experiments, the gains to this sort of matching asymptote to zero, so these techniques will primarily benefit smaller advertisers.

The future is also looking up for evaluating television advertising and associated “cross media” interaction effects. More people are viewing TV through devices like the Xbox and through services like Google TV, both of which link users to ads in systems similar to major web publishers. Furthermore, these users often have identifiers such as Google or Microsoft usernames that can link television, sponsored search, and display ads for a single individual. Never before in the history of advertising has this been possible. The ability to measure cross-channel effects with the reliability of randomized experiments opens the door to many new questions for academics and many new strategies for advertisers. As more forms of advertising become measurable on an individual level, our ability to provide reliable estimates of advertising effectiveness will expand as well. The advances so far have already set a new state-of-the-art in measurement, and we expect the trend to continue.

## 9 Conclusion

The science of measuring advertising effectiveness has evolved considerably due to new digital data sources and experimentation platforms. Experiments using online media have not only allowed us to cleanly (albeit often noisily) measure the impact of advertising, but also laid bare how difficult adfx really is to measure. Compared to effects often studied in applied economics, such as the returns to schooling, price elasticity and the impact of subsidies/taxes, the effects of advertising explain many orders of magnitude less of the variation in the dependent variable. For individual-level sales impact, even a *successful* campaign will typically have an  $R^2$  lower than  $0.00001$ . Accordingly, selection effects that explain very little variation can swamp treatment effects; this means observational methods are not suited for measuring adfx—we advocate strongly for experiments. In Section 5

we give a specific example of one such bias we call “activity bias,” which is a temporal selection bias that fells observational methods attempting to account for heterogeneity between treated and untreated groups.

We view experimentation on the individual level with the ad delivery linked to purchasing behavior as a true game-changer offered by digital media as compared to traditional counterparts. Whether in search or display, new experimental platforms can give feedback to advertisers that is immune from the biases that plague observational methods. Another important advance is computational advertising. Computational advertising helps solve the targeting problem and usefully shifts the conversation from “who to hit” to “what do I get.” Yet neither of these advances solve all the measurement problems in advertising science. Experiments are noisy and computational advertising relies on humans to enter the key parameters, such as valuations of clicks or conversions, that govern spend.

Moving forward, experimentation and data collection technology is evolving alongside new forms of ad-serving. Questions such as the cross-derivative of certain media on the effectiveness of other media (“cross-channel effects”) will be in play in the coming years. Measuring the effectiveness of media, such as television, that were previously not technologically feasible, because randomizing delivery was not possible at scale, will also greatly expand knowledge on advertising effectiveness. This will in turn allow firms to more accurately guide their advertising expenditure. But all is not rosy: we have shown too that certain questions such as the long-run effects of advertising and the impact of “brand advertising” appear to be out of reach for at least the next 5–10 years. We await new developments in data collection and experimentation technology to facilitate the answers to these and new questions.

## References

- Abraham, M. (2008). The off-line impact of online ads. *Harvard Business Review*, 86(4):28.
- Bagwell, K. (2005). The economic analysis of advertising. *Handbook of Industrial Organization*, 3.
- Blake, T., Nosko, C., and Tadelis, S. (2013). Consumer Heterogeneity and Paid Search Effectiveness: A Large Scale Field Experiment. In *NBER Working Paper*, pages 1–26. NBER.
- BLS (2010). In *American Time Use Survey*, volume <http://www.bls.gov/tus/charts/leisure.htm>.
- Broder, A. (2008). Computational advertising and recommender systems. In *Proceedings of the 2008 ACM conference on Recommender systems*, pages 1–2. ACM.
- Card, D. (1999). The causal effect of education on earnings. *Handbook of labor economics*, 3:1801–1863.

- Cary, M., Das, A., Edelman, B., Giotis, I., Heimerl, K., Karlin, A., Mathieu, C., and Schwarz, M. (2007). Greedy bidding strategies for keyword auctions. In *Proceedings of the 8th ACM conference on Electronic commerce*, pages 262–271. ACM.
- Chan, D., Ge, R., Gershony, O., Hesterberg, T., and Lambert, D. (2010). Evaluating online ad campaigns in a pipeline: causal models at scale. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 7–16. ACM.
- Dickie, M., Fisher, A., and Gerking, S. (1987). Market transactions and hypothetical demand data: A comparative study. *Journal of the American Statistical Association*, pages 69–75.
- Goldfarb, A. and Tucker, C. (2011). Search engine advertising: Channel substitution when pricing ads to context. *Management Science*, 57(3):458–470.
- Johnson, G., Lewis, R. A., and Reiley, D. H. (2012). Location, location, location: Geo-targeting increases effectiveness of online display advertising. *Unpublished manuscript*.
- Johnson, G. H., Lewis, R. A., and Reiley, D. H. (2013). Add more ads? experimentally measuring incremental purchases due to increased frequency of online display advertising. In *Working paper*.
- Kumar, D. and Yildiz, T. (2011). Measuring online ad effectiveness. In *12th ACM Conference on Electronic Commerce*. ACM.
- Lewis, R. (2010). *Where’s the “Wear-Out?”: Online Display Ads and the Impact of Frequency*. PhD thesis, MIT PhD Dissertation.
- Lewis, R., Rao, J., and Reiley, D. (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. and Reiley, D. (2013a). Super bowl advertising causes down-to-the-minute online search behavior. In *Proceedings of the 14th ACM Conference on Electronic Commerce*. ACM.
- Lewis, R., Reiley, D., and Schreiner, T. (2012). Ad attributes and attribution: Large-scale field experiments measure online customer acquisition. *Unpublished manuscript*.
- Lewis, R. A. and Nguyen, D. T. (2013). A samsung ad and the ipad: Display advertising’s spillovers to online search. *Unpublished manuscript*.
- Lewis, R. A. and Rao, J. M. (2013). On the near-impossibility of measuring the returns to advertising. *Unpublished manuscript*.

- Lewis, R. A. and Reiley, D. H. (2012). Advertising effectively influences older users: A yahoo! experiment measuring retail sales. *Review of Industrial Organization*, (forthcoming).
- Lewis, R. A. and Reiley, D. H. (2013b). Online advertising and offline sales: Measuring the effects of retail advertising via a controlled experiment on yahoo! *Unpublished manuscript*.
- 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.
- Mincer, J. (1962). On-the-job training: Costs, returns, and some implications. *The Journal of Political Economy*, 70(5):50–79.
- Murphy, J., Allen, P., Stevens, T., and Weatherhead, D. (2005). A meta-analysis of hypothetical bias in stated preference valuation. *Environmental and Resource Economics*, 30(3):313–325.
- Pandey, S., Chakrabarti, D., and Agarwal, D. (2007). Multi-armed bandit problems with dependent arms. In *Proceedings of the 24th International Conference on Machine learning*, pages 721–728. ACM.
- Rusmevichientong, P. and Williamson, D. (2006). An adaptive algorithm for selecting profitable keywords for search-based advertising services. In *Proceedings of the 7th ACM Conference on Electronic commerce*, pages 260–269. ACM.
- Tucker, C. (2012). Social advertising. *Available at SSRN 1975897*.