Observational studies with unknown time of treatment

Guillaume W. Basse[†] Alexander Volfovsky[†] Edoardo M. Airoldi*

^{*}Guillaume W. Basse is a graduate student in the Department of Statistics at Harvard University (gbasse@fas.harvard.edu). Alexander Volfovsky is an National Science Foundation Postdoctoral Fellow in the Department of Statistics at Harvard University (volfovsky@fas.harvard.edu). Edoardo M. Airoldi is an Associate Professor of Statistics at Harvard University (airoldi@fas.harvard.edu). This work was partially supported by the National Science Foundation under grants CAREER IIS-1149662, IIS-1409177, and DMS-1402235, and by the Office of Naval Research under grant YIP N00014-14-1-0485. Edoardo M. Airoldi is an Alfred P. Sloan Research Fellow, and a Shutzer Fellow at the Radcliffe Institute for Advanced Studies. The authors wish to thank Michael Els, Robert Haslinger, Mark Lowe, Sean Murphy, and Donald B. Rubin for providing data, comments, and valuable insights.

Abstract

Time plays a fundamental role in causal analyses, where the goal is to quantify the effect of a specific treatment on future outcomes. In a randomized experiment, times of treatment, and when outcomes are observed, are typically well defined. In an observational study, treatment time marks the point from which pre-treatment variables must be regarded as outcomes, and it is often straightforward to establish. Motivated by a natural experiment in online marketing, we consider a situation where useful conceptualizations of the experiment behind an observational study of interest lead to uncertainty in the determination of times at which individual treatments take place. Of interest is the causal effect of heavy snowfall in several parts of the country on daily measures of online searches for batteries, and then purchases. The data available give information on actual snowfall, whereas the natural treatment is the anticipation of heavy snowfall, which is not observed. In this article, we introduce formal assumptions and inference methodology centered around a novel notion of plausible time of treatment. These methods allow us to explicitly bound the last plausible time of treatment in observational studies with unknown times of treatment, and ultimately yield valid causal estimates in such situations.

Keywords: Causal inference; Rubin causal model; Plausible time of treatment; Natural experiment; Online marketing.

Contents

1	Inti	roduction	1
2	A n	notion of plausible time of treatment	2
	2.1	Formal assumptions	4
3	Des	signing an observational study with unknown time of treatment	6
	3.1	Conceptualization of treatment	6
	3.2	Pilot study to identify plausible times of treatment	7
	3.3	Main study	9
4	Sim	ulations	10
	4.1	Design choices	10
		4.1.1 Non-negative effect curve, zero at observation time	11
		4.1.2 Non-negative effect curve, positive at observation time	12
		4.1.3 Positive-negative effect curve, zero at observation time	12
	4.2	Analysis of the results	13
5	Analyzing the effect of snowfall on online behavior		15
	5.1	The data	16
	5.2	Sensitivity analysis	18
	5.3	A heuristic for the choice of α and ϵ	21
6	Cor	ncluding remarks	23
\mathbf{R}_{i}	References		

1 Introduction

Observational studies are common across many disciplines, including the health and social sciences where practitioners are interested in assessing the causal effect of a non-randomized treatment (Rosenbaum, 2002, 2010). For example, a drug company might solicit information on the blood pressure of individuals following the voluntary ingestion of a particular drug, or a fast food restaurant might record observations on sales over time. The former might be interested in the effect of the drug on blood pressure, while the latter might be interested in the effect of the forecast of an adverse weather event on sales. In both cases, the data are gathered from historical records, thus without the ability to randomize treatment. Because of the lack of randomization, naive estimators of the average treatment effect that consider the difference between the average effect of treated and untreated individuals are likely biased for the effect of interest (Rubin, 1991; Imbens and Rubin, 2015).

To make the study of causal effects concrete an analyst is required to define the treatment of interest, the outcome of interest that she believes the treatment might affect, times at which to conceptualize treatment happened, and times at which to measure the outcomes (Splawa-Neyman et al., 1990; Rubin, 1974). Within the context of a randomized experiment the meaning and definition of these three components are very straightforward, but in the context of an observational study much greater care is required in defining them to allow for proper causal inference. In particular, one might only observe a proxy for treatment and so the exact treatment and, more importantly, the exact time of treatment are unknown. When this is the case it becomes unclear when the outcome of interest should have been measured in order to make a proper causal statement. This is the setting we consider in this paper.

As a concrete example, we consider studying the relationship between the adverse weather in February and March of 2015 across the United States and online battery searches at a major US retailer. The relationship between weather and sales has been studied in the lit-

erature (e.g., see Murray et al., 2010; Starr-McCluer, 2000; Zwebner et al., 2013), with the sometimes implicit assumption of a clearly defined treatment (e.g., temperature, exposure to sunlight) and so causal statements are reserved for post-weather measurements. In our case we are interested in studying whether the perception of future extreme snow increases conversions but we do not observe this perception. Having observed the snow the methods proposed in Section 3 determines the time points prior to the snow event for which causal statements can be made. In a novel data set, provided by a large online advertising technology firm, we have measurements of these outcomes across February and March of 2015 for 79 designated marketing areas (DMA) across the United States. For each of these super-metropolitan areas we have covariate information pertaining to the demographics of the people living in the area, as well as measurement of snow accumulation for each day.

The rest of the paper is organized as follows: In Section 2 we briefly outline the formal notation of the Rubin Causal Model and we introduce two new assumptions, which complement the assumptions employed in causal inference, that enable causal analyses of data that are missing an exact time of treatment. Section 3 describes the proposed three stage methodology. Section 4 provides simulation results. Details of the marketing data set and the data analysis are reported in Section 5. Remarks follow, in Section 6.

2 A notion of plausible time of treatment

Throughout this paper we restrict ourselves to two levels of treatment and for convenience we refer to "control" for the lower level of treatment, and to "treatment" for the higher. As in the classical causal inference literature Z_i denotes treatment indicator for an individual i and is equal to 1 if the unit is assigned to treatment at time T_i^{treat} and is 0 otherwise (Imbens and Rubin, 2015). Potential outcomes for each unit are denoted by $Y_{i,T_i^{\text{treat}},t}(Z_i=0)$ or $Y_{i,T_i^{\text{treat}},t}(Z_i=1)$ for control and treated levels of treatment when treatment occurs at time

 T_i^{treat} and the outcomes are observed at time $t > T_i^{\text{treat}}$.

In contrast with the classical setting, the treatment indicator Z_i and the time of treatment T_i^{treat} remain hidden and instead, D_i and T_i^{event} are observed. D_i is the indicator that an event associated with the treatment occurred at time T_i^{event} . In particular it is not necessary that $Z_i = D_i$ and it is likely that $T_i^{\text{treat}} < T_i^{\text{event}}$. That is, the event that is observed is not necessarily a perfect proxy for the actual treatment and the actual treatment time occurs before the event. In what follows we formulate a set of assumptions on the relationship between Z_i , D_i , T_i^{treat} and T_i^{event} that allow us to determine the times t for which an observed quantity $Y_{i,t}^{\text{obs}}$ corresponds to the potential outcome under treatment or under control:

$$Y_{i,t}^{\text{obs}} = D_i \cdot Y_{i,T_i^{\text{treat}},t}(Z_i = 1) + (1 - D_i) \cdot Y_{i,T_i^{\text{treat}},t}(Z_i = 0)$$
(1)

We then employ these assumptions to reconstruct an observational data set that is suited for causal inference.

To develop these assumptions we require a new notion that we refer to as the "Last Plausible Randomized Experiment Time" (termed LaPRET and represented by T_i^{LaPRET}). In the context of an idealized randomized experiment, outlined in the left hand panel of Figure 1, LaPRET defines the first time point at which the two potential outcomes are differentiable. In other words, in a randomized experiment where the treatment has an effect after being administered at treatment time T_i^{treat} , if the outcome were to be observed at time $T_i^{\text{obs}} \leq T_i^{\text{LaPRET}}$ then no difference between $Y_{i,T_i^{\text{treat}},T_i^{\text{obs}}}(Z_i=0)$ and $Y_{i,T_i^{\text{treat}},T_i^{\text{obs}}}(Z_i=1)$ would be discernible (that is, a treatment effect would be undetectable). On the other hand observing at time $T_i^{\text{obs}} > T_i^{\text{LaPRET}}$ would yield a non-zero treatment effect. It is clear that when there is no treatment effect, T_i^{LaPRET} is not well defined – the method proposed in this paper for identifying T_i^{LaPRET} in observational studies accounts for this issue.

2.1 Formal assumptions

We require two new assumptions to discuss Equation (1). The first facilitates the comparison between the realized treatment and the associated event while the second provides a rationale for identifying the times at which treatment could have happened.

Assumption 1. (Unit Treatment Status Identification Strength) For treatment indicator Z_i and associated event indicator D_i assume that

$$cor(Z_i, D_i) \ge 1 - \eta \ \forall i \text{ for } \eta \text{ small.}$$

Assumption 2. (Constant Unit Time of Treatment) For LaPRET, T_i^{LaPRET} , treatment time T_i^{treat} and associated event time T_i^{event} assume

(i)
$$d_i = T_i^{\text{event}} - T_i^{\text{LaPRET}} = d > 0$$
 for all i

(ii)
$$T_i^{\text{treat}} \leq T_i^{\text{LaPRET}}$$
 for all i

The statement of Assumption 1 reflects potential uncertainty about realized treatment status based on event status. In the extreme setting where $\eta = 0$ the assumption states that $Z_i = D_i$ for all units i, fully revealing the treatment indicator. This is the setting when a longitudinal study reveals information about a treatment that occurred during a previous wave of the study. In the less extreme setting of η small, the assumption provides a generative method for multiple imputed data sets all of which are suitable for causal inference (see e.g. Rubin (1996)). In particular, if we reconceptualize Equation 1 as:

$$\Pr\left(Y_{i,t}^{\text{obs}} = D_i, Y_{i,T_i^{\text{treat}},t}(Z_i = 1) + (1 - D_i)Y_{i,T_i^{\text{treat}},t}(Z_i = 0)\right) \ge 1 - \delta \tag{2}$$

we have $\delta = \delta(\eta)$ which informs the interpretation of a sensitivity analysis that varies η . When $\eta \neq 0$ there are multiple versions of the "complete data" as there is now a nonzero probability that $D_i = 1$ and $Z_i = 0$. This condition can easily be accommodated by constructing multiple data sets in which the correlation between Z_i and D_i is $1-\eta$ and reporting causal estimates across these data sets (Rubin, 1996).

Assumption 2 is an identification assumption for LaPRET that explicitly addresses the lack of information about the time of treatment. The first part of the assumption states that the time between an individual LaPRET and the time of the associated event is equal for all individuals in the population. As such, by observing two similar individuals, one in the control and one in the treatment groups (using the information from Assumption 1), one can compute the LaPRET as illustrated in the right hand panel of Figure 1. The second part of the assumption states that the time of treatment is always strictly prior to LaPRET. meaning that observations at the LaPRET are potential outcomes and so can be used to compute causal estimates. Our methodology relies on the first part of this assumption - we estimate \hat{d} from a pilot study and then leverage this information in the larger population.

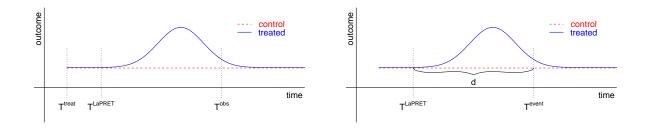


Figure 1: This is an illustration of Assumption 2. The left panel hand describes the scenario of an idealized experiment – T_i^{treat} is the known treatment time, T_i^{LaPRET} is the LaPRET and T_i^{obs} is the observation time. The right hand panel describes the observational data set where an event associated with treatment is observed at time T_i^{event} and d (of Assumption 2), the time between the event time and the LaPRET is the same for all i (and hence does not require a subscript).

If this assumption is violated then the estimate of \hat{d} is unreliable and care must be taken to insure estimates are causal. In such situations, illustrated in the simulation study, a smaller value of d can be chosen to insure that the estimate is causal. This suggests a relaxation of the condition in part (i) of Assumption 2, to the get the following.

Assumption 2'. (Stable Unit Time of Treatment) For LaPRET, T_i^{LaPRET} , treatment time T_i^{treat} and associated event time T_i^{event} assume

(i)
$$d_i = T_i^{\text{event}} - T_i^{\text{LaPRET}} > 0$$
 for all i

(ii)
$$T_i^{\text{treat}} \leq T_i^{\text{LaPRET}}$$
 for all i

In this more general case, there exists d > 0 such that $d_i = T_i^{\text{obs}} - T_i^{\text{LaPRET}} > d$ for all i.

3 Designing an observational study with unknown time of treatment

In this section, we describe how one may leverage the assumptions of Section 2 in order to take a large data set that does not include treatment indicators or time of treatment indicators and to produce a potentially reduced data set that can be used for conceptualizing and performing a complete observational study (Rosenbaum, 2002, 2010).

3.1 Conceptualization of treatment

Unlike in classical causal inference where treatment and treatment time are known to the scientist, in this setting a proper notion of treatment must be conceptualized. The issue here is similar to the one faced in the perceived treatment literature where one cannot use race or sex as the treatment in an experiment but instead "perceived race" or "perceived sex" can

be used (Greiner and Rubin, 2011). We also rely on the notion of "perceived treatment". For example, in the application to the advertising data in Section 5 the treatment is the expectation of an extreme snow event happening in the future perceived by individuals. The notion of perception is needed here as the true treatment is never observed in our setting and so we reconstruct it based on an event that would have been foreshadowed by such a perception. This is exactly the type of correlation between the treatment indicator Z_i and the event indicator D_i that Assumption 1 describes quantitatively. It must be noted that D_i identifies both the treatment and the control groups and so it is likely that some units will be discarded in order to better satisfy Assumption 1. This data reduction step is explored in detail in Section 5 where the only units allowed to be considered as control units are those that always experience less than a thresholded amount of snow.

Once we have identified the event that serves as proxy for the true treatment we must identify the correlation in Assumption 1. In a longitudinal study setting, where a later wave question (serving as the indicator D_i) might ask if an individual received an intervention at a previous wave (the treatment indicator Z_i) the correlation can be set to one. In more ambiguous situations, such as the one discussed in the Section 4, a sensitivity analysis based on this correlation should be performed.

3.2 Pilot study to identify plausible times of treatment

Once the treatment is defined and a variable D_i is identified we need to find the individual LaPRET values T_i^{LaPRET} . Since these are only identifiable from the joint distribution of outcomes we must infer those from the data. However, if we use a data driven approach that considers the complete data set we would be violating a fundamental principle of causal inference that does not allow parameters in the analysis to depend on the post-treatment data.

To overcome this difficulty we propose to perform a pilot study to identify T_i^{LaPRET} . An example of this approach was recently undertaken by Wager and Athey (2015) to apply decision tree methods to causal inference. Under the pilot study we consider a small sample of individuals that are identified as treated and as control units. Once this sample is chosen we match treated and control pairs. The choice of a particular matching mechanism depends on the applied problem and the available observed covariates associated with each unit—in our simulations and applied example we employ propensity score matching (Rosenbaum and Rubin, 1983). In a slight abuse of notation, the matched pairs are now identified by the subscript i. Leveraging Assumption 2(i) (or 2(i)') we can identify the difference between the time of the event D_i and the LaPRET, $d_i = T_i^{\text{event}} - T_i^{\text{LaPRET}}$, as it is assumed to be constant (or greater than a non-zero constant) for all units. As such finding d_i becomes a problem of identifying time points t where $Y_{i,t}^{\text{obs}}|D_i = 1$ and $Y_{i,t}^{\text{obs}}|D_i = 0$ are close. Letting $\Delta_{i,t}$ equal that difference for matched pair i at time t we say that the LaPRET for matched pair i is

$$\hat{T}_{i}^{\text{LaPRET}} \in \underset{t < T_{i}^{\text{event}}}{\text{arg max}} \left\{ t \text{ s.t. } |\Delta_{i,t}| < \max_{t} \frac{|\Delta_{i,t}|}{\alpha} \text{ and} \right.$$

$$\text{s.t. } |\partial \Delta_{i,t}| < \epsilon, \ |\partial \Delta_{i,s_{1}}| > \epsilon, \ |\partial \Delta_{i,s_{2}}| < \epsilon \text{ for some } s_{2} > s_{1} > t \right\}.$$

That is, the estimated LaPRET for pair i is at the maximal time point such that the difference $\Delta_{i,t}$ is smaller than a $(1/\alpha)$ fraction of the maximal difference and the rate of change of $\Delta_{i,t}$ with respect to time is small but there exists a later time point where the rate of change is large. The parameter α captures the expected variability in the values of $Y_{i,t}^{\text{obs}}$. Small values of α allow for larger differences between treated and control observations to be evaluated as no treatment effect. As such, larger values of α lead to more conservative $\hat{T}_i^{\text{LaPRET}}$ – those that are closer to T_i^{event} . The parameter ϵ controls the rate of change of $\Delta_{i,t}$. It insures that the procedure is able to differentiate between no effect (where $Y_{i,t}(1) = Y_{i,t}(0) \ \forall t$) and a

situation where the effect of treatment either induces volatility or decreases to zero before the time of observation. Small values of ϵ require the volatility between the treated and control observations to be small almost all the time while larger values allow for lots of volatility but require that a bigger volatility event occurs. As such both too small and too large values of ϵ lead to extremely conservative behavior. We study the behavior of α and ϵ in a simulation study.

After computing $\hat{T}_i^{\text{LaPRET}}$ we can construct \hat{d} for Assumption 2 using some function of the set $\{T_i^{\text{event}} - \hat{T}_i^{\text{LaPRET}}\}_i$. In the simulation study we explore using the average of those differences to choose \hat{d} . Choosing the minimum forms a conservative estimate of d that accommodates the relaxed Assumption 2'.

In practice, the pilot study can be performed on a subset of units treated at a period of time before or after the main study happened.

3.3 Main study

The final step of the pipeline is the construction of a possibly reduced data set from the units that were not used in the pilot study which allows for valid causal inference (Rosenbaum, 2002, 2010; Imbens and Rubin, 2015). In particular, having identified the relationship between the treatment indicator Z_i and the event indicator D_i as well as the latency between LaPRET and the event in the previous two steps one constructs the data set as follows: (1) among the remaining units, construct a matched sample then (2) for each matched pair, discard information that is recorded prior to $\hat{T}_i^{\text{LaPRET}} = T_i^{\text{event}} - \hat{d}$. The remaining data then represents the potential outcomes under treatment and control.

4 Simulations

The purpose of our method is to give the analyst an interval of days \hat{d} prior to the event, for which he can make causal statements about the conceptualized treatment. This section has two objectives: first, we assess how the choice of parameters α and ϵ affects the behavior of our method in different scenarios. Second, we explore how this behavior is affected by different types and levels of noise in the observations.

4.1 Design choices

We conduct three experiments of increasing complexity, trying capturing different real life response profiles. For each of the three scenarios, we consider a pair of "idealized response surfaces" under control and treatment – that is, response surfaces with no noise that capture one of the three scenarios of interest. We denote pair of responses corresponding to the k^{th} scenario by $(\mu_0^{(k)}(t), \mu_1^{(k)}(t))$. For all three scenarios, the idealized control response surface is the flat line at zero, $\mu_0^{(k)}(t) = 0$. The idealized response surfaces are represented in Figure 2.

We also consider two potential sources of noise in the observations. First we simulate the fact the potential outcomes are only noisy versions of the idealized response surface. So for each scenario k, we generate N treatment and control potential outcomes surfaces:

$$Y_{i,t}^{(k)}(a) \sim \text{Normal } (\mu_a^{(k)}(t), \sigma^2)$$
 (3)

for a=0,1, and i=1...N. The noise parameter σ was given the following values $\sigma=0.005,0.01,0.015,0.02$, but such that the potential outcome surfaces generated in a given simulation all shared the same parameter σ . The second source of noise we consider is with the time of observation T_i^{event} of the event, which we simulate with a contamination model: $T_i^{\text{event}} = T_i^{\text{LaPRET}} + d + c$ where T_i^{LaPRET} and d are fixed and c is a contamination model that

is governed by one of the following distributions:

$$f_1(c) = 0 \text{ w.p. } 1$$

$$f_2(c) = \begin{cases} -1 & \text{w.p. } 0.25 \\ 0 & \text{w.p. } 0.5 \\ 1 & \text{w.p. } 0.25 \end{cases}$$

$$f_3(c) = \begin{cases} -1 & \text{w.p. } 0.1 \\ 0 & \text{w.p. } 0.5 \\ 1 & \text{w.p. } 0.4 \end{cases}$$

$$f_4(c) = \begin{cases} -1 & \text{w.p. } 0.4 \\ 0 & \text{w.p. } 0.5 \\ 1 & \text{w.p. } 0.1 \end{cases}$$

For each of the $4 \times 4 = 16$ possible noise structures (σ^2, f_i) , we study the behavior of our method for $\alpha = 1, 6, \dots 96$ and $\epsilon = 0.0001, 0.02, 0.2, 0.3, 0.4, 0.5$. So for each combination of parameters $(\sigma^2, f_i, \alpha, \epsilon)$, we simulated a data set with N = 600 observations, and compute \hat{d} as in step~2 of Section 3. We forgo matching in the simulations since this is beyond the scope of our contributions. The three simulations are described in more details below, and the results are analyzed in Section 4.2.

4.1.1 Non-negative effect curve, zero at observation time

The first simulation, we consider the scenario in the left panel of Figure 2, where the LaPRET is $T_i^{\text{LaPRET}} = 3$, the event is observed at time $T_i^{\text{event}} = 14$ (and so $d_i = T_i^{\text{event}} - T_i^{\text{LaPRET}} = 11$) and the treatment response surface is represented, for mathematical convenience, by the function $\mu_1^{(1)}(t) = \max \left\{ 0, \sin(\frac{2\pi}{15}(t-3.5)) \right\}$. In this setup, the true effect is always nonnegative, but is zero at the observation time, and so a naive analysis might find no effect even when one is present. This version of a response is likely when the natural event is catastrophic—if individuals expect a snow storm, for instance, they might order batteries and snow tires in advance so as to have them ready for after the storm.

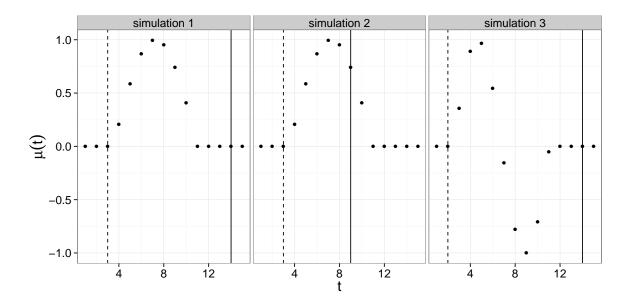


Figure 2: Response surface for the three simulations. The solid lines are the observation times, T_i^{event} , and the dotted lines are the LaPRET, T_i^{LaPRET}

4.1.2 Non-negative effect curve, positive at observation time

The second simulation considers the same response surface as Simulation 1 but a shifted observation time $T_i^{\text{event}} = 9$ such that d = 6 as displayed in the second panel of Figure 2. In this scenario, the true effect is still always non-negative, but the observation time corresponds to a point where there is still a difference between the treated and control levels of the potential outcomes. This is a possible response surface for insurance quote requests due to an upcoming storm—there is an increase due to the forecast and it does not necessarily go back down to the previous levels until after the storm. Similar response surfaces were observed in for online marketing campaigns (Lewis et al., 2011).

4.1.3 Positive-negative effect curve, zero at observation time

The third simulation introduces volatility into the response surface. Here $\mu_1^{(k)}(t)$ is as in (4), $T_i^{\text{LaPRET}} = 2$ and $T_i^{\text{event}} = 14$. In this scenario, the true effect is positive immediately after

 T_i^{LaPRET} , but then changes sign before reaching zero shortly before T_i^{event} . In particular, there is a point of zero effect between T_i^{LaPRET} and T_i^{event} which does not correspond to the LaPRET. In the context of medical trials, this volatility could correspond to the side-effect of a drug on a person's blood pressure. It could also correspond to the purchasing of commodities before a big storm – the dip illustrated in the right panel of Figure 2 corresponds to the fact that once an individual stocks up on a commodity, he is likely to buy less of it for a while. Another interpretation is that in anticipation of the event, individuals move their usual purchasing day to before the event, provoking a dip in the days immediately preceding the event.

$$u(t) = \begin{cases} 0 & \text{if } v(t) \le 0 \text{ and } t < 4 \text{ or if } v(t) \ge 0 \text{ and } t > 10\\ \sin(\frac{3.5\pi}{15}(t - 2.5)) & \text{otherwise.} \end{cases}$$
(4)

4.2 Analysis of the results

The effect of the different parameters is similar in all three scenarios, so we only describe the results of the second Simulation, described in Section 4.1.2. The results are summarized in Figure 3. Figures 4 and 5 summarize the results of the first and third simulations, respectively. We start by noticing that although the contamination models have a small impact on the aggregated value of \hat{d} , this can have a large impact on the integer part of \hat{d} . This quantity, which we denote by $\lfloor \hat{d} \rfloor$, is the quantity of interest since it represents the number of days before the event date for which we can make causal statements. For instance in Figure 3, for low levels of α , and $\sigma = 0.005$, we see that $\lfloor \hat{d} \rfloor = 6$ for contamination models f_1, f_2 and f_3 , but $\lfloor \hat{d} \rfloor = 5$ for f_4 . This uncertainty is difficult to account for, and our method is very sensitive to it. For fixed values of σ and ϵ , we see that \hat{d} is close to the true value d for low values of α , but then decreases as α increases. For very low values of α , however, the \hat{d} decreases. For fixed values of α and σ , \hat{d} is close to zero for low values of ϵ , increases

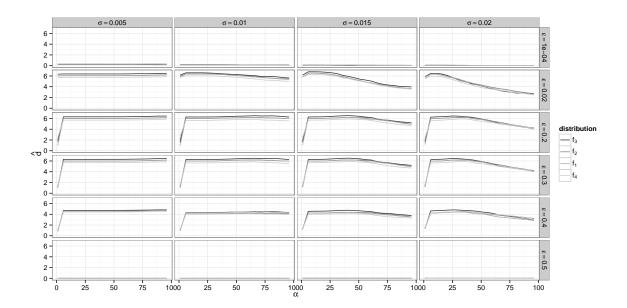


Figure 3: Estimate \hat{d} as a function α under different contamination models, different levels of noise, and different values of ϵ for simulation 2, described in Section 4.1.2.

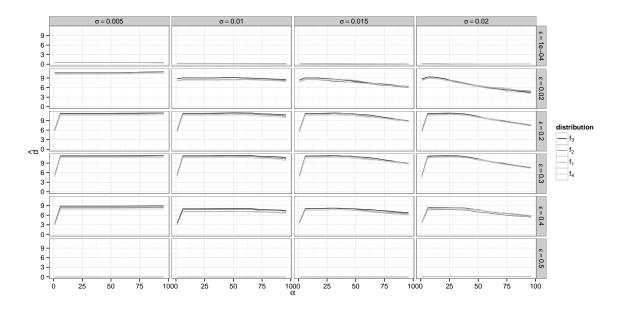


Figure 4: Estimate \hat{d} as a function α under different contamination models, different levels of noise, and different values of ϵ for simulation 1, described in Section 4.1.1.

until a certain point with ϵ , then decreases again to reach zero for high values of ϵ . Finally, we see that as the noise level σ increases, \hat{d} decreases, especially for high values of α .

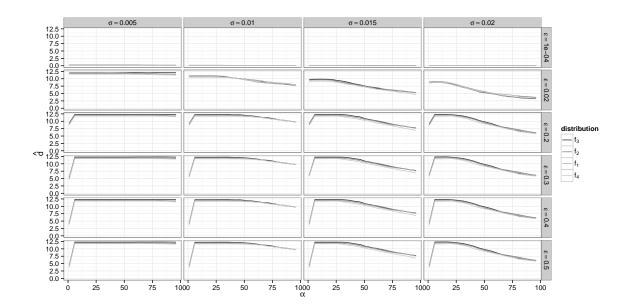


Figure 5: Estimate \hat{d} as a function α under different contamination models, different levels of noise, and different values of ϵ for simulation 3, described in Section 4.1.3.

A decrease in the value of \hat{d} is a conservative behavior from a causal perspective, as it reduces the range of effects that we can call causal. This means, based on our observations in the previous paragraph, that as the noise σ increases, the estimator becomes increasingly conservative in these scenarios, which is a desirable behavior. This property of our estimator is further explored in Section 5. Our observations have also made clear the fact that increasing the value of the parameters α tends to make the estimator more conservative.

5 Analyzing the effect of snowfall on online behavior

In this section we estimate the causal treatment effect of perceived large quantities of future snow on sales of products on the internet. Our estimates are based on data provided by MaxPoint Interactive Inc., an advertising technology company based in Raleigh, North Carolina. The data was provided according to designated marketing area (DMA) which corresponds to 79 super-metropolitan areas in the United States. For each DMA, the data

contains three types of information over a period of two months: (1) searches for batteries on a major retailer's website, which we will consider as the outcome of our analysis, (2) Demographic information from the census bureau, summarized in Figure 6 and Figure 7, and (3) the cumulative daily snowfalls. This data set is interesting as it can be seen as recording a natural experiment at the nationwide scale, in which the treatment time is unknown (e.g., see Angrist et al., 2000; Dunning, 2012; Phan and Airoldi, 2015).

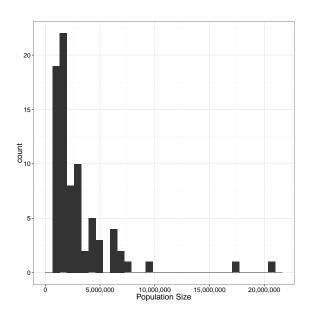


Figure 6: Distribution of the population of the 79 DMAs considered.

5.1 The data

Due to the large populations and geographical extents that different DMAs cover, it is unreasonable to believe that an aggregate of weather at the DMA level will constitute a treatment for all units in said DMA. Another way to say this is that DMAs are not homogenous when it comes to weather, and there can be a significant weather event happening in a DMA that will affect only a small portion of the units inside that DMA. We thus build a synthetic scenario that better illustrates the methodology we propose. We introduce the concept of tradezones which can be thought of as smaller subergions of DMAs, within which weather

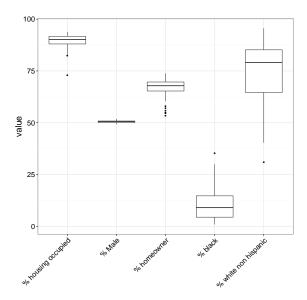


Figure 7: Distribution of the demographic variables across the 79 DMAs considered will be more homogeneous. For each DMA, a number of tradezones proportional to its population, totaling 3676 synthetic tradezones over all DMAs.

For tradezone j in DMA i, we create for each day t a synthetic observation from a Normal $(\mu_i(t), \sigma^2)$ for $y_{ij}(t)$ truncated below at zero, where $\mu_i(t)$ is the real observed outcome for DMA i on day t, such that $y_{ij}(t) \geq 0$ for all t. Varying the noise level σ allows us to evaluate the robustness of our procedure. We ran simulations for $\sigma = 2^k$, k = 1...7. To give some perspective about the relative size of the noise, we report that the median outcome in the data across all DMAs and all days is 13, and the 95th percentile is 60.

Another element required to construct a realistic data set is to handle the weather, keeping in mind that the original motivation for introducing the tradezones was to deal with weather heterogeneity within DMAs. To address this, we set a threshold $h = 1 \,\mathrm{kg/m^2}$ of snow, and for each DMA, we considered the days for which the snow precipitations exceeded that threshold. Suppose there where K such days for DMA i, which we will label $\{t_1, \ldots, t_K\}$, and let S_1, \ldots, S_K be the corresponding precipitations in $\mathrm{kg/m^2}$. For each tradezone j, we selected a day of observation $T^{(\mathrm{event},j)}$ at random among $\{t_1, \ldots, t_K\}$, such that $p(T^{(\mathrm{event},j)})$

 t_l) = $\frac{S_l}{\sum_k S_k}$, that is, proportional to snow precipitations in the eligible days. Tradezones in DMAs for which no snow precipitation exceeded $l = 0.3 \,\mathrm{kg/m^2}$ were all assigned to the control group, while all remaining tradezones were ignored.

With this definition of treatment, it is reasonable to say that in most cases, no unit with $D_i = 0$ would have had the perception that it is going to be hit by an event such as $1 \,\mathrm{kg/m^2}$ of snow. This definition justifies setting the correlation in Assumption 1 to be one (that is $\eta = 0$) and hence the observed outcomes are realizations of the potential outcomes (that is $\delta = 0$ in Eq (2)). Larger h strengthen this assumption but drastically reduce the number of DMAs in treatment.

Census data at the DMA level are used to complete the synthetic data set. For each DMA i, we have a vector of covariates x_i . Since the propensity score matching is performed at the tradezone level, we let tradezone j in DMA i inherit its covariate vector $x_i^{(j)}$ from the parent DMA. That is, we assume that $x_i^{(j)} = x_i \,\forall j$.

In summary, for each DMA (for which we have real data), we have generated synthetic tradezones, for which we simulated synthetic observations, and selected a day of observation for the snow event based on the distribution of snowfall in the DMA. Our purpose is to illustrate how an analyst would apply our method to that kind of data, and what kind of robustness he should be expecting. The results we report below represent causal estimates of nationwide battery searches online (for a major american retailer) based on the synthetic data.

5.2 Sensitivity analysis

For each of the seven levels of noise $\sigma = 2^k$, k = 1...7, we generated one data set as described in Section 5.1, and ran both a pilot study and an analysis as described in Section 3, with parameters $\alpha = 2.5$ and $\epsilon = 4$. Injecting different levels of noise assesses the sensitivity of

the analysis to the synthetic data generating process. We discuss heuristics for the choice of α and ϵ in Section 5.3.

To stabilize the data and account for difference in baseline outcomes among DMAs, we consider lagged differences between outcomes. That is, if $y_{ij}(t)$ is the outcome of tradezone j in DMA i at time t, then we considered the transformed outcomes $y_{ij}^*(t) = y_{ij}(t) - y_{ij}(t-1)$ (we ignore the first day in the data set), and applied the method in Section 5.1 to the transformed data. This analysis provides insight into the causal changes in behavior from day to day due to perceived future snow. Figure 8 shows the values of \hat{d} obtained in the seven pilot studies, where each pilot study contains 878 of the 3676 tradezones. The solid line shows the integer part of \hat{d} which is the real quantity of interest, since d is the number of days before the event for which we can report the effect as being causal. We see that the \hat{d} (and hence the solid line) decreases as the standard deviation of the noise increases in our simulations. This confirms the conservative behavior identified in Section 4: the noisier the data, the more conservative we become about making causal statements.

After obtaining values for \hat{d} from the pilot studies, we completed the analysis for each

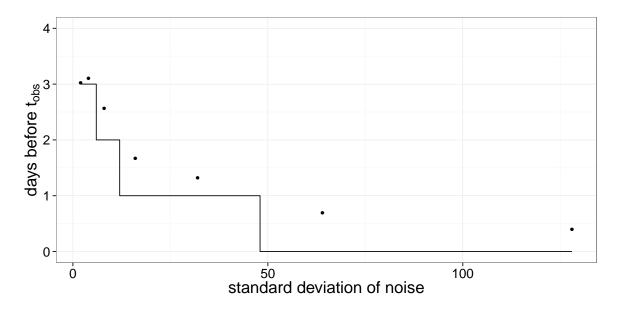


Figure 8: Sensitivity of \hat{d} to noise in the tradezones

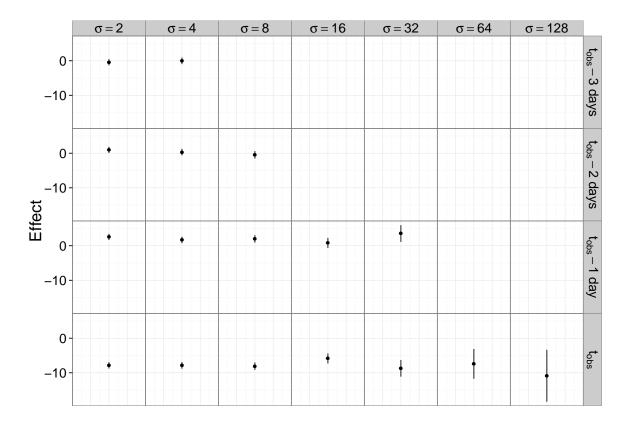


Figure 9: Sensitivity of the estimated effect to noise in the tradezones. Each frame represents the average causal effect and a 95% confidence interval.

data set based on the remaining 2798 tradezones. Figure 9 summarizes the results, and illustrate the conservative nature of our procedure: for high levels of noise, the only effect that can be reported as causal is that on the day of the observed event. For low values of the noise however, causal statements can be made up to two days prior to the day of the weather event.

From this data we conclude that there is a causal relationship between searches for batteries and perceived future snow events that takes on the form of the ATE in Simulation 3 (see Fig 2). That is, the causal effect appears to be consistently positive and growing several days prior to the weather event (at all noise levels) and then decreases drastically on the day of the event. For example, at the lowest noise level introduced into the synthetic data

the day-to-day change is close to zero three days prior to the event, but increases to 0.95 two days before the event and again increases to 2.55 one day before the event. From the day before the event to the day of the event there is a drop of 7.8 in the rate of searches on average.

5.3 A heuristic for the choice of α and ϵ

We have seen in Section 4 that the choice of α and ϵ governs how conservative our method is. Although it is ultimately up to the analyst to chose and justify the parameters he uses in the analysis, we provide some heuristics to guide this choice. In this section, we will let $Y_{i,t}^{obs}$ be either the observed outcomes, or the first order differences which we denoted by Y^* in the previous section. Consider $\Delta_{i,t}$ and $\partial \Delta_i$, t as in Section 3, then let $\bar{\Delta}$ and $\bar{\partial}\bar{\Delta}$ be the respective averages of their absolute values, $\tilde{\Delta}$ and $\bar{\partial}\bar{\Delta}$ the respective maxima of their absolute values, and se(Δ) and se(Δ) the respective standard errors of their absolute values. We suggest choosing values of α and ϵ satisfying:

$$\alpha_{min} = \frac{\tilde{\Delta}}{\bar{\Delta} + 3se(\Delta)} \le \alpha \le \frac{\tilde{\Delta}}{\bar{\Delta} + se(\Delta)} = \alpha_{max}$$
 (5)

and

$$\epsilon_{min} = \frac{\widetilde{\partial \Delta}}{\overline{\partial \Delta} + 3\operatorname{se}(\partial \Delta)} \le \epsilon \le \frac{\widetilde{\partial \Delta}}{\overline{\partial \Delta} + \operatorname{se}(\partial \Delta)} = \epsilon_{max}$$
 (6)

These heuristics are based on the interpretation of α and ϵ as measures of variation in the outcomes and the first differences of the outcomes. The choice of α_{min} is the ratio of the maximum absolute variation in outcomes to three standard deviations more than the mean while α_{max} is the ratio of the maximum variation in outcomes to one standard deviation more than the mean. The smaller α values are thus associated with how extreme the maximum is in comparison to the mean. Similarly for ϵ , smaller values are associated with how extreme the maximum of the absolute value of first differences in outcomes is in comparison to the

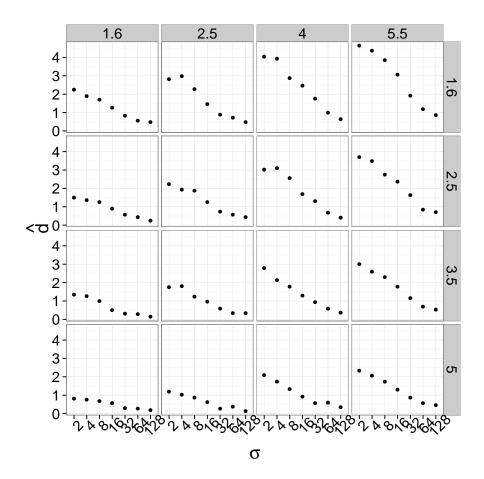


Figure 10: Values of \hat{d} for a range of parameters α and ϵ obtained using the heuristics in Section 5.3, for different levels of noise

mean. The larger α and ϵ values have similar interpretation but with respect to the less extreme single standard deviation from the mean. The choice of one and three standard deviations is motivated by normal asymptotics.

With our synthetic data, the ranges of $[\alpha_{min}, \alpha_{max}]$ and $[\epsilon_{min}, \epsilon_{max}]$ depend on the noise level σ — for illustration purposes, we consider the ranges $\alpha \in [1.6, 5.5]$ and $\epsilon \in [1.6, 5]$ obtained by the unions of the ranges for the different values $\sigma = 2^i$, i = 1...7. Figure 10 displays the values of \hat{d} that would be obtained for different combinations of α and ϵ within the range of our heuristics. Note that the parameters α and ϵ do not affect the estimate of the effect, only their causal interpretation.

6 Concluding remarks

The methodology we developed in this article is not intended to supersede any of the traditional methodology for dealing with observational studies, but rather to complement it. At a very high level, one can see our method as a pre-processing step, which provide the analyst with one level of protection against unsubstantiated causal claims.

We make two main contributions. First we provide a set of assumptions and a method to determine a window before the day of the observed event for which we can we can make causal statement. Our second contribution is cleanly separate the overall process into a pilot study, which is used to determine the window in which we can make causal statements, and the causal analysis, which is carried on a disjoint subset of data. This precaution insulates the causal analysis from any dependence on the observed outcomes used into the analysis. We have shown in simulation studies that our method becomes increasingly conservative when the observed outcomes become volatile, and that passed a certain level of noise, the method precludes any causal statement beyond the date when the treatment proxy is observed.

Our methodology extends the reach of causal inference to a specific type of observational studies, in which it is suspected that the causal effect happens before the date in which a proxy to the treatment is observed. The price paid for this extension is a reliance on extra assumptions, and a loss of efficiency since the causal analysis is carried only on a subset of data. We also emphasize the fact that the causal analysis carried being an observational studies, it suffers from the usual limitations.

References

J. D. Angrist, K Graddy, and G. W. Imbens. The interpretation of instrumental variables estimators in simultaneous equations models with an application to the demand for fish.

- Review of Economics Studies, 67(3):499–527, 2000.
- Thad Dunning. Natural Experiments in the Social Sciences: A Design-Based Approach.

 Cambridge University Press, Cambridge; New York, October 2012. ISBN 9781107698000.
- D James Greiner and Donald B Rubin. Causal effects of perceived immutable characteristics.

 Review of Economics and Statistics, 93(3):775–785, 2011.
- Guido Imbens and Donald B. Rubin. Causal inference for statistics, social, and biomedical sciences: an introduction. Cambridge University Press, New York, 2015. ISBN 978-0521885881.
- Randall A Lewis, Justin M Rao, and David H Reiley. 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, 2011.
- Kyle B Murray, Fabrizio Di Muro, Adam Finn, and Peter Popkowski Leszczyc. The effect of weather on consumer spending. *Journal of Retailing and Consumer Services*, 17(6): 512–520, 2010.
- T. Q. Phan and E. M. Airoldi. A natural experiment of social network formation and dynamics. *Proceedings of the National Academy of Sciences*, 112(21):6595–6600, 2015.
- P. R. Rosenbaum. Observational Studies. Springer, 2nd edition, 2002.
- P. R. Rosenbaum. Design of Observational Studies. Springer, 2nd edition, 2010.
- Paul R Rosenbaum and Donald B Rubin. Assessing sensitivity to an unobserved binary covariate in an observational study with binary outcome. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 212–218, 1983.
- D. B. Rubin. Practical implications of modes of statistical inference for causal effects and the critical role of the assignment mechanism. *Biometrics*, 47:1213–1234, 1991.

- Donald B Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology*, 66(5):688, 1974.
- Donald B Rubin. Multiple imputation after 18+ years. Journal of the American statistical Association, 91(434):473–489, 1996.
- Jerzy Splawa-Neyman, DM Dabrowska, TP Speed, et al. On the application of probability theory to agricultural experiments. essay on principles. section 9. *Statistical Science*, 5(4): 465–472, 1990.
- Martha Starr-McCluer. The effects of weather on retail sales. Divisions of Research & Statistics and Monetary Affairs, Federal Reserve Board, 2000.
- Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using random forests. arXiv preprint arXiv:1510.04342, 2015.
- Yonat Zwebner, Leonard Lee, and Jacob Goldenberg. The temperature premium: Warm temperatures increase product valuation. *Journal of Consumer Psychology*, 24(2):251–259, 2013.