

CAUSAL INFERENCE WITH SPATIO-TEMPORAL DATA:

ESTIMATING THE EFFECTS OF AIRSTRIKES ON INSURGENT VIOLENCE IN IRAQ*

Georgia Papadogeorgou[†] Kosuke Imai[‡] Jason Lyall[§] Fan Li[¶]

The first version: April 27, 2020

This version: July 16, 2021

Abstract

Many causal processes have spatial and temporal dimensions. Yet the classic causal inference framework is not directly applicable when the treatment and outcome variables are generated by spatio-temporal processes with an infinite number of possible event locations. We extend the potential outcomes framework to these settings by formulating the treatment point process as a stochastic intervention. Our causal estimands include the expected number of outcome events in a specified area under a particular stochastic treatment assignment strategy. We develop methodology that allows for arbitrary patterns of spatial spillover and temporal carryover effects. Using martingale theory, we show that the proposed estimator is consistent and asymptotically normal as the number of time periods increases, even when the propensity score is estimated. We propose a sensitivity analysis for the possible existence of unmeasured confounders, and extend it to the Hájek estimator. Simulation studies are conducted to examine the estimators' finite sample performance. Finally, we use the proposed methods to estimate the effects of American airstrikes on insurgent violence in Iraq from February 2007 to July 2008. We find that increasing the average number of daily airstrikes for up to one month results in more insurgent attacks across Iraq and within Baghdad. We also find evidence that airstrikes can displace attacks from Baghdad to new locations up to 400 kilometers away.

Keywords: carryover effects, inverse probability of treatment weighting, point process, sensitivity analysis, spillover effects, stochastic intervention, unstructured interference

*Lyall gratefully acknowledges financial support from the Air Force Office of Scientific Research (Grant #FA9550-14-1-0072). The findings and conclusions reached here do not reflect the official views or policy of the United States Government or Air Force. In addition, Imai thanks the Sloan Foundation (# 2020-13946) for financial support. The authors would also like to thank Soubhik Barari, Iavor Bojinov, Naoki Egami, Connor Jerzak, Sayar Karmakar, and Neil Shephard for their constructive comments.

[†]Assistant Professor, Department of Statistics, University of Florida, Gainesville FL 32611. Email: gpapadogeorgou@ufl.edu, URL: <https://gpapadogeorgou.netlify.com>

[‡]Professor, Department of Government and Department of Statistics, Harvard University. 1737 Cambridge Street, Institute for Quantitative Social Science, Cambridge MA, 02138. Email: imai@Harvard.Edu, URL: <https://imai.fas.harvard.edu>

[§]James Wright Chair in Transnational Studies and Associate Professor, Department of Government, Dartmouth College, Hanover, NH 03755. Email: jason.lyall@dartmouth.edu, URL: www.jasonlyall.com

[¶]Professor, Department of Statistical Science, Duke University, Durham, NC 27708. Email: fl35@duke.edu, URL: <http://www2.stat.duke.edu/~fl35>

1 Introduction

Many causal processes involve both spatial and temporal dimensions. Examples include the environmental impact of newly constructed factories, the economic and social effects of refugee flows, and the various consequences of disease outbreaks. These applications also illustrate key methodological challenges. First, when the treatment and outcome variables are generated by spatio-temporal processes, there exists an infinite number of possible treatment and event locations at each point in time. In addition, spatial spillover and temporal carryover effects are likely to be complex and may not be well understood.

Unfortunately, the classical causal inference framework that dates back to [Neyman \(1923\)](#) and [Fisher \(1935\)](#) is not directly applicable to such settings. Indeed, standard causal inference approaches assume that the number of units that can receive the treatment is finite (e.g., [Rubin, 1974](#); [Robins, 1997](#)). Although a small number of studies develop a continuous time causal inference framework, they do not incorporate a spatial dimension (e.g., [Gill and Robins, 2001](#); [Zhang *et al.*, 2011](#)). In addition, causal inference methods have been used for analyzing functional magnetic resonance imaging (fMRI) data, which have both spatial and temporal dimensions. For example, [Luo *et al.* \(2012\)](#) apply randomization-based inference, while [Sobel and Lindquist \(2014\)](#) employ structural modelling. We instead focus on data generated by different underlying processes, leading to new estimands and estimation strategies.

Specifically, we consider settings in which the treatment and outcome events are assumed to be generated by spatio-temporal point processes (Section 3). The proposed method is based on a single time series of spatial patterns of treatment and outcome variables, and builds upon three strands of the causal inference literature: interference, stochastic interventions, and time series.

First, we address the possibility that treatments might affect outcomes at a future time period and at different locations in arbitrary ways. Although some researchers have considered unstructured interference, they assume non-spatial and cross-sectional settings (see e.g., [Basse and Airoidi, 2018](#); [Sävje *et al.*, 2019](#), and references therein). In addition, [Aronow *et al.* \(2019\)](#) study spatial randomized exper-

iments in a cross-sectional setting, and under the assumption that the number of potential intervention locations is finite and their spatial coordinates are known and fixed. By contrast, our proposed spatio-temporal causal inference framework allows for *temporally and spatially unstructured interference* over an infinite number of locations.

Second, instead of separately estimating the causal effects of treatment received at each location, we consider the impacts of different *stochastic treatment assignment strategies*, defined formally as the intervention distributions over treatment point patterns. Stochastic interventions have been used to estimate effects of realistic treatment assignment strategies (Díaz Muñoz and van der Laan, 2012; Young *et al.*, 2014; Papadogeorgou *et al.*, 2019) and to address challenging causal inference problems including violation of the positivity assumption (Kennedy, 2019), interference (Hudgens and Halloran, 2008; Imai *et al.*, 2021), mediation analysis (Lok, 2016; Díaz and Hejazi, 2019), and multiple treatments (Imai and Jiang, 2019). We show that this approach is also useful for causal inference with spatio-temporal treatments and outcomes.

Finally, our methodology allows for arbitrary patterns of spatial and temporal interference. As such, our estimation method does not require the separation of units into minimally interacting sets (e.g., Tchetgen Tchetgen *et al.*, 2017). Nor does it rely on an outcome modelling approach that entails specifying a functional form of spillover effects based on, for example, geographic distance. Instead, we view our data as a single time series of maps, which record the locations of treatment and outcome realizations as well as the geographic coordinates of other relevant events. Our estimation builds on the time-series causal inference approach pioneered by Bojinov and Shephard (2019).

We propose a spatially-smoothed inverse probability weighting estimator that is consistent and asymptotically normal under a set of reasonable assumptions, regardless of whether the propensity scores are known, or estimated from a correctly specified model (Section 4). To do so, we establish a new central limit theorem for martingales that can be widely used for causal inference in observational, time series settings. We also show that the proposed estimator based on the estimated propensity score has a lower

asymptotic variance than when the true propensity score is known. This generalizes the existing theoretical result under the independently and identically distributed setting (Hirano *et al.*, 2003) to the spatially and temporally dependent setting. Finally, to assess the potential impact of unobserved confounding, we develop a sensitivity analysis method by generalizing the sensitivity analysis of Rosenbaum (2002) to our spatio-temporal context and to the Hájek estimator with standardized weights (Section 5). We conduct simulation studies to assess the finite sample performance of the proposed estimators (Section 6).

Our motivating application is the evaluation of the effects of American airstrikes on insurgent violence in Iraq from February 2007 to July 2008 (Section 2). We consider all airstrikes during each day anywhere in Iraq as a *treatment pattern*. Instead of focusing on the causal effects of each airstrike, we estimate the effects of different *airstrike strategies*, defined formally as the distributions of airstrikes throughout Iraq (Section 7). The proposed methodology enables us to capture spatio-temporal variations in treatment effects, shedding new light on how airstrikes affect the location, distribution, and intensity of insurgent violence.

Specifically, we show that a higher number of airstrikes, without modifying their spatial distribution, increases the number of insurgent attacks, especially near Baghdad, Mosul, and the roads between them. We also find that changing the focal point of airstrikes to Baghdad without modifying the overall frequency shifts insurgent attacks from Baghdad to Mosul and its environs. These findings demonstrate that airstrikes can increase insurgent attacks *and* disperse them over considerable distances. Furthermore, our analysis suggests that increasing the number of airstrikes may initially reduce attacks but ultimately increases them over the long run. Our sensitivity analysis indicates, however, that these findings are somewhat sensitive to the potential existence of unmeasured confounders.

The proposed methodology has a wide range of applications beyond the specific example analyzed in this paper. For example, the causal effects of pandemics and crime on a host of economic and social outcomes could be evaluated using our methodology. With the advent of massive and granular data sets, we expect the need to conduct causal analysis of spatio-temporal data will only continue to grow.

2 Motivating Application: Airstrikes and Insurgent Activities in Iraq

Airstrikes have emerged as a principal tool for fighting against insurgent and terrorist organizations in civil wars around the globe. In the past decade alone, the United States has conducted sustained air campaigns in at least six different countries, including Afghanistan, Iraq, and Syria. Although it has been shown that civilians have all-too-often borne the brunt of these airstrikes (Lyall, 2019b), we have few rigorous studies that evaluate the impact of airstrikes on subsequent insurgent violence. Even these studies have largely reached opposite conclusions, with some claiming that airpower reduces insurgent attacks while others arguing they spark escalatory spirals of increased violence (e.g., Lyall, 2019a; Mir and Moore, 2019a; Dell and Querubin, 2018; Kocher *et al.*, 2011).

Moreover, all existing studies have two interrelated methodological shortcomings: they carve continuous geographic space into discrete, often arbitrary, units, and they make simplifying assumptions about patterns of spatial and temporal interference. Mir and Moore (2019b), for example, argue that drone strikes in Pakistan have reduced terrorist violence. But they use a coarse estimation strategy that bins average effects of drone strikes into broad half-year increments over entire districts that cannot capture local spatial and temporal dynamics. Similarly, Rigterink (2021) draws on 443 drone strikes to estimate airstrike effects on 13 terrorist groups in Pakistan, concluding that they have mixed effects. Yet her group-month estimation strategy cannot detect spillover effects nor accurately capture the timing of insurgent responses. In short, we need a flexible methodological approach that avoids the pitfalls of binning treatment and outcome measures into too-aggregate, possibly misleading, temporal and spatial units.

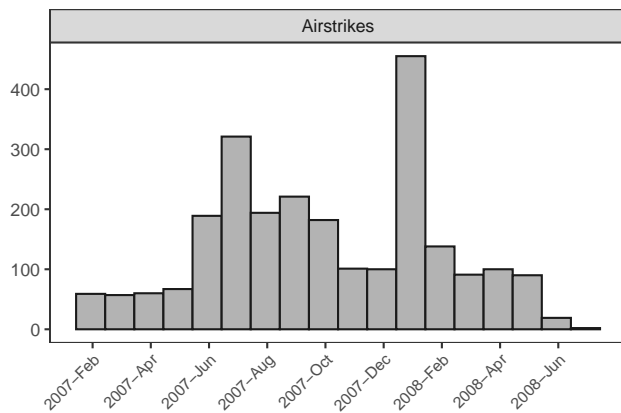
We enter this debate by examining the American air campaign in Iraq. We use declassified US Air Force data on airstrikes and shows of force (simulated airstrikes where no weapons are released) for the February 2007 to July 2008 period. The period in question coincides with the “surge” of American forces and airpower designed to destroy multiple Sunni and Shia insurgent organizations in a bid to turn

the war's tide.

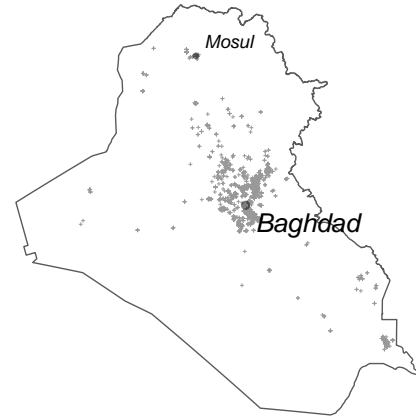
Aircraft were assigned to bomb targets via two channels. First, airstrikes were authorized in response to American forces coming under insurgent attack. These close air support missions represented the vast majority of airstrikes in 2007–08. Second, a small percentage (about 5%) of airstrikes were pre-planned against high-value targets, typically insurgent commanders, whose presence had been detected from intercepted communications or human intelligence. In each case, airstrikes were driven by insurgent attacks that were either ongoing or had occurred in the recent past in a given location. As a result, the models used later in this paper adjust for prior patterns of insurgent violence in a given location for several short-term windows.

We also account for prior air operations, including shows of force, by American and allied aircraft. Insurgent violence in Iraq is also driven by settlement patterns and transportation networks. Our models therefore include population size and location of Iraqi villages and cities as well as proximity to road networks, where the majority of insurgent attacks were conducted against American convoys. Finally, prior reconstruction spending might also drive the location of airstrikes. Aid is often provided in tandem with airstrikes to drive out insurgents, while these same insurgents often attack aid locations to derail American hearts-and-minds strategies. Taken together, these four factors—recent insurgent attacks, the presence of American forces, settlement patterns, and prior aid spending—drove decisions about the location and severity of airstrikes.

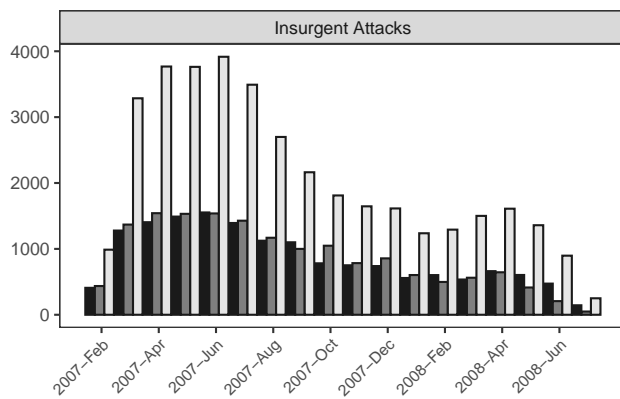
Figure 1 summarizes the spatial and temporal distributions of airstrikes (treatment variable) and insurgent violence (outcome variable). Figure 1a presents the temporal distribution of airstrikes recorded by the US Air Force each month. There were a total of 2,246 airstrikes during this period. Figure 1b plots the spatial density of these airstrikes across Iraq, with spatial clustering observed around Baghdad and the neighboring “Sunni Triangle,” a hotspot of insurgency. Figure 1c plots the monthly distribution of insurgent attacks by type: Improvised Explosive Devices (IEDs), small arms fire (SAF), and other attacks. A total of 68,573 insurgent attacks were recorded by the US Army's CIDNE database during



(a) Airstrikes over time



(b) Airstrikes over space



(c) Insurgent violence over time



(d) Insurgent violence over space

Figure 1: Distribution of the treatment and outcome point processes over time and space. Plots (b) and (d) show the locations of airstrikes and insurgent attacks, respectively, during the time period February 23, 2007 to July 05, 2008. Insurgent attacks are sorted into one of three categories: Improvised Explosive Devices (IEDs), Small Arms Fire (SAF), and other attacks.

this time period. Finally, Figure 1d plots the locations of insurgent attacks across Iraq. Baghdad, the Sunni Triangle, and the highway leading north to Mosul are all starkly illustrated.

3 Causal Inference Framework for Spatio-temporal Data

In this section, we propose a causal inference framework for spatio-temporal point processes. We describe the setup, and define causal estimands based on stochastic interventions.

3.1 The Setup

We represent the locations of airstrikes for each time period (e.g., day) as a spatial point pattern measured at time $t \in \mathcal{T} = \{1, 2, \dots, T\}$ where T is the total number of the discrete time periods. Let $W_t(s)$ denote the binary treatment variable at location s for time period t , indicating whether or not the location receives the treatment during the time period. We use W_t as a shorthand for $W_t(\Omega)$, which evaluates the binary treatment variable $W_t(s)$ for each element s in a set Ω . The set Ω is *not* assumed to be a finite grid, but it is allowed to include an infinite number of locations that may receive the treatment. In addition, \mathcal{W} represents the set of all possible point patterns at each time period where, for simplicity, we assume that this set does not vary across time periods, i.e., $W_t \in \mathcal{W}$ for each t . The set of *treatment-active locations*, i.e., the locations that receive the treatment, at time t is denoted by $S_{W_t} = \{s \in \Omega : W_t(s) = 1\}$. We assume that the number of treatment-active locations is finite for each time period, i.e., $|S_{W_t}| < \infty$ for any t . In our study, the treatment-active locations correspond to the set of coordinates of airstrikes. Finally, $\overline{W}_t = (W_1, W_2, \dots, W_t)$ denotes the collection of treatments over the time periods $1, 2, \dots, t$.

We use w_t to represent a realization of W_t and $\overline{w}_t = (w_1, w_2, \dots, w_t)$ to denote the history of treatment point pattern realizations from time 1 through time t . Let $Y_t(\overline{w}_t)$ represent the potential outcome at time $t \in \mathcal{T}$ for any given treatment sequence $\overline{w}_t \in \mathcal{W}^t = \mathcal{W} \times \dots \times \mathcal{W}$, depending on *all* previous treatments. Similar to the treatment, $Y_t(\overline{w}_t)$ represents a point pattern with locations $S_{Y_t(\overline{w}_t)}$, which are referred to as the *outcome-active locations*. In our study, $S_{Y_t(\overline{w}_t)}$ represents the locations of insurgent attacks if the patterns of airstrikes had been \overline{w}_t . Let $\overline{\mathcal{Y}}_T = \{Y_t(\overline{w}_t) : \overline{w}_t \in \mathcal{W}^t, t \in \mathcal{T}\}$ denote the collection of potential outcomes for all time periods and for all treatment sequences.

Among all of these potential outcomes for time t , we only observe the one corresponding to the

observed treatment sequence, denoted by $Y_t = Y_t(\overline{\mathbf{W}}_t)$. We use $\overline{\mathbf{Y}}_t = \{Y_1, Y_2, \dots, Y_t\}$ to represent the collection of observed outcomes up to and including time period t . In addition, let \mathbf{X}_t be the set of possibly time-varying confounders that are realized prior to W_t but after W_{t-1} . No assumption is necessary about the temporal ordering of any variables in \mathbf{X}_t and Y_{t-1} . Let $\overline{\mathcal{X}}_T = \{\mathbf{X}_t(\overline{\mathbf{w}}_{t-1}) : \overline{\mathbf{w}}_{t-1} \in \mathcal{W}^{t-1}, t \in \mathcal{T}\}$ be the set of potential values of \mathbf{X} under any possible treatment history and for all time periods. We also assume that the observed covariates correspond to the covariates under the observed treatment path, $\mathbf{X}_t = \mathbf{X}_t(\overline{\mathbf{W}}_{t-1})$, and use $\overline{\mathbf{X}}_t = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_t)$ to denote the collection of observed covariates over the time periods $1, 2, \dots, t$. Finally, we use $\overline{H}_t = \{\overline{\mathbf{W}}_t, \overline{\mathbf{Y}}_t, \overline{\mathbf{X}}_{t+1}\}$ to denote all observed history preceding the treatment at time $t + 1$.

Since our statistical inference is based on a single time series, we consider all potential outcomes and potential values of the time-varying confounders as *fixed*, pre-treatment quantities. Then, the randomness we quantify is with respect to the assignment of treatment W_t given the complete history including all counterfactual values \overline{H}_{t-1}^* where $\overline{H}_t^* = \{\overline{\mathbf{W}}_t, \overline{\mathcal{Y}}_T, \overline{\mathcal{X}}_T\}$ and $\overline{H}_t \subset \overline{H}_t^*$.

3.2 Causal Estimands under Stochastic Interventions

A notion central to our proposed causal inference framework is *stochastic intervention*. Instead of setting a treatment variable to a fixed value, a stochastic intervention specifies the probability distribution that generates the treatment under a potentially counterfactual scenario. Although our framework accommodates a large class of intervention distributions, for concreteness, we consider intervention distributions based on Poisson point processes, which are fully characterized by an intensity function $h : \Omega \rightarrow [0, \infty)$. For example, a homogeneous Poisson point process with $h(s) = h$ for all $s \in \Omega$, implies that the number of treatment-active locations follows a $\text{Poisson}(h|\Omega|)$ distribution, with locations distributed independently and uniformly over Ω . In general, the specification of stochastic intervention should be motivated by policy or scientific objectives. Such examples in the context of our study are given in Section 7.1.

Our causal estimands are the expected number of (potential) outcome-active locations under a specific stochastic intervention of interest, and the comparison of such quantities under different intervention

distributions. We begin by defining the causal estimands for a stochastic intervention taking place over a single time period. Let F_h denote the distribution of a spatial point process with intensity h . Also, let $N_B(\cdot)$ denote a counting measure on a region $B \subset \Omega$. Then, we can define the expected number of outcome-active locations for a region B at time t as

$$N_{Bt}(F_h) = \int_{\mathcal{W}} N_B(Y_t(\overline{\mathbf{W}}_{t-1}, w_t)) dF_h(w_t) = \int_{\mathcal{W}} \left| S_{Y_t}(\overline{\mathbf{W}}_{t-1}, w_t) \cap B \right| dF_h(w_t). \quad (1)$$

In our application, this quantity represents the expected number of insurgent attacks within a region of Iraq B if the airstrikes at time t were to follow the point process specified by F_h , given the observed history of airstrikes up to time $t - 1$. The region B does not need to be defined as a connected subset of Ω , and it can be the union of potentially non-bordering sets (for example, the suburbs of two cities).

We can extend the above estimand to an intervention taking place over M consecutive time periods. Consider an intervention, denoted by $F_{\mathbf{h}} = F_{h_1} \times \cdots \times F_{h_M}$, under which the treatment at time t is assigned according to F_{h_1} , at time $t - 1$ according to F_{h_2} , continuing until time period $t - M + 1$ for which treatment is assigned according to F_{h_M} . A treatment path based on this intervention is displayed in Figure 2(a). Then, we define a general estimand as

$$\begin{aligned} N_{Bt}(F_{\mathbf{h}}) &= \int_{\mathcal{W}^M} N_B(Y_t(\overline{\mathbf{W}}_{t-M}, w_{t-M+1}, \dots, w_t)) dF_{h_1}(w_t) \cdots dF_{h_M}(w_{t-M+1}) \\ &= \int_{\mathcal{W}^M} \left| S_{Y_t}(\overline{\mathbf{W}}_{t-M}, w_{t-M+1}, \dots, w_t) \cap B \right| dF_{h_1}(w_t) \cdots dF_{h_M}(w_{t-M+1}). \end{aligned} \quad (2)$$

This quantity represents the expected number of outcome events within region B and at time t if the treatment point pattern during the previous M time periods was to follow the stochastic intervention with distribution $F_{\mathbf{h}}$. Treatments during the initial $t - M$ time periods were the same as observed. A special case of $F_{\mathbf{h}}$ assumes that treatments during the M time periods are independent and identically distributed draws from the same distribution F_h , which we denote by $F_{\mathbf{h}} = F_h^M$.

Given the above setup, we define the average treatment effect of stochastic intervention $F_{\mathbf{h}'}$ versus

$F_{\mathbf{h}''}$ for a region B at time t as

$$\tau_{Bt}(F_{\mathbf{h}'}, F_{\mathbf{h}''}) = N_{Bt}(F_{\mathbf{h}''}) - N_{Bt}(F_{\mathbf{h}'}), \quad (3)$$

where $\mathbf{h}' = (h'_1, h'_2, \dots, h'_M)$ represents a collection of treatment intensities over M consecutive time periods (similarly for \mathbf{h}'').

We further consider the average, over time periods $t = M, M + 1, \dots, T$, of the expected potential outcome for region B at each time period if treatments during the M preceding time periods arose from $F_{\mathbf{h}}$. This quantity is defined as

$$N_B(F_{\mathbf{h}}) = \frac{1}{T - M + 1} \sum_{t=M}^T N_{Bt}(F_{\mathbf{h}}). \quad (4)$$

Figure 2 shows two of the terms averaged in Equation (4), i.e., $N_{Bt}(F_{\mathbf{h}})$ and $N_{B(t+1)}(F_{\mathbf{h}})$. For $N_{Bt}(F_{\mathbf{h}})$, treatments up to $t - M$ are set to their observed values, and treatments at time periods $t - M + 1, \dots, t$ are drawn from $F_{\mathbf{h}}$. The same definition applies to $N_{B(t+1)}(F_{\mathbf{h}})$, but intervention time periods are shifted by 1: treatments up to $t - M + 1$ are set to their observed values, while treatments during time periods

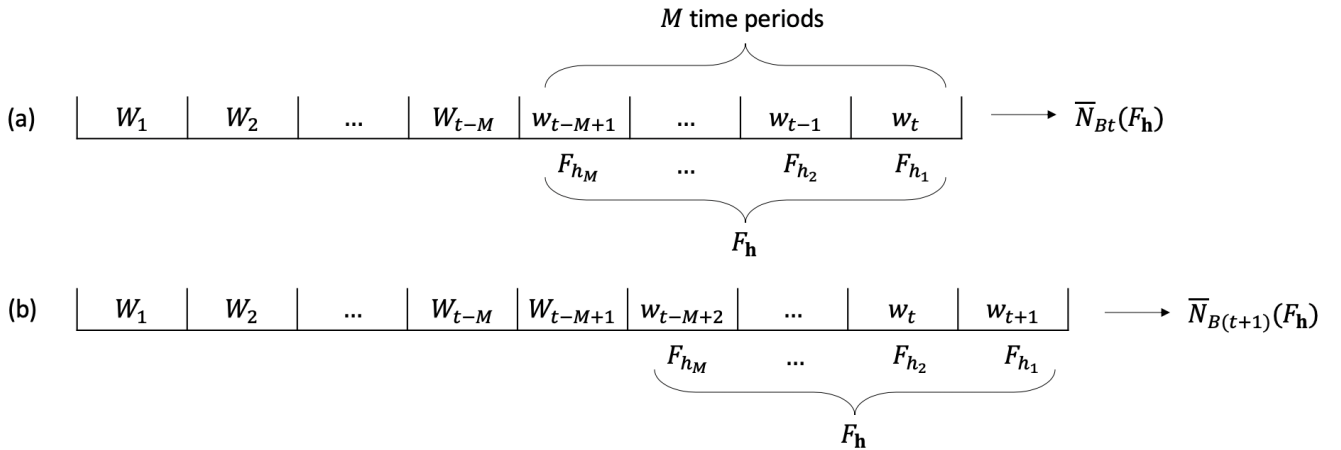


Figure 2: Graphical Illustration of Stochastic Intervention over Multiple Time Periods for Time Period t and $t + 1$. Under intervention $F_{\mathbf{h}}$, treatments during time periods $t - M + 1, \dots, t - 1, t$ are assigned according to distributions $F_{h_M}, \dots, F_{h_2}, F_{h_1}$.

$t - M + 2, \dots, t + 1$ are drawn from $F_{\mathbf{h}}$. In Equation (4), the summation starts at $t = M$ since the quantity $N_{Bt}(F_{\mathbf{h}})$ assumes that there exist M prior time periods during which treatments are intervened on. We suppress the dependence of $N_B(F_{\mathbf{h}})$ on T for notational simplicity.

Similarly, based on $N_B(F_{\mathbf{h}})$, we define the causal effect of intervention $F_{\mathbf{h}'}$ versus $F_{\mathbf{h}''}$ as

$$\tau_B(F_{\mathbf{h}'}, F_{\mathbf{h}''}) = N_B(F_{\mathbf{h}''}) - N_B(F_{\mathbf{h}'}) = \frac{1}{T - M + 1} \sum_{t=M}^T \tau_{Bt}(F_{\mathbf{h}'}, F_{\mathbf{h}''}). \quad (5)$$

This estimand represents the average, over time periods $t = M, M + 1, \dots, T$, of the expected change in the number of points at each time period when the observed treatment path $\overline{\mathbf{W}}_T$ was followed until $t - M$ with subsequent treatments W_{t-M+1}, \dots, W_t arising according to $F_{\mathbf{h}'}$ versus $F_{\mathbf{h}''}$.

The effect size of a point pattern treatment would depend on M , and a greater value of M allows one to study slow-responding outcome processes. Moreover, specifying $F_{\mathbf{h}'}$ and $F_{\mathbf{h}''}$ such that they are identical except for the assignment at M time periods prior, h'_M, h''_M , yields the lagged effect of a treatment change, which resembles the lagged effects defined by [Bojinov and Shephard \(2019\)](#) for binary treatments and non-stochastic interventions.

The above estimands are defined while conditioning on the treatments of all previous time periods. This is important because we do not want to restrict the range of temporal carryover effects. Although the proposed estimand is generally data-dependent, the quantity becomes fixed under some settings. For example, if the potential outcomes at time t are restricted to depend at most on the latest L treatment point patterns, then the estimands for stochastic interventions that take place over $M \geq L$ time periods will no longer depend on the observed treatment path.

4 Estimation and Inference

In this section, we introduce a set of causal assumptions and the proposed estimator that combines inverse probability of treatment weighting with kernel smoothing. We then derive its asymptotic properties. All proofs are given in Appendix B.

4.1 The Assumptions

Similar to the standard causal inference settings, variants of the unconfoundedness and overlap assumptions based on stochastic interventions are required for the proposed methodology. For simplicity, we focus on stochastic interventions with identical and independent distribution over M periods, $F_{\mathbf{h}} = F_h^M$, and intensity h . Our theoretical results, however, extend straightforwardly to stochastic interventions with non-i.i.d. treatment patterns.

Assumption 1 (Unconfoundedness). *The treatment assignment at time t does not depend on any, past or future, potential outcomes and potential confounders conditional on the observed history of treatments, confounders and outcomes up to time $t - 1$:*

$$f(W_t | \overline{\mathbf{W}}_{t-1}, \overline{\mathcal{Y}}_T, \overline{\mathcal{X}}_T) = f(W_t | \overline{H}_{t-1}).$$

Assumption 1 resembles the sequential ignorability assumption in the standard longitudinal settings (Robins, 1999; Robins *et al.*, 2000), but it is more restrictive. The assumption requires that the treatment assignment does not depend on both past and future potential values of the time-varying confounders as well as those of the outcome variable, conditional on their past observed values. In contrast, the standard sequential ignorability assumption only involves future potential outcomes.

Unfortunately, sequential ignorability would not suffice in the current setting. The reason is that we utilize data from a *single* unit measured repeatedly over many time periods to draw causal conclusions. This contrasts with the typical longitudinal settings where data are available on a large number of independent units over a short time period. Our assumption is similar to the non-anticipating treatment assumption of Bojinov and Shephard (2019) for binary non-stochastic treatments, while explicitly showing the dependence on the time-varying confounders. By requiring the treatment to be conditionally independent of the time-varying confounders, we assume that all “back-door paths” from treatment to either the outcome or the time-varying confounders are blocked (Pearl, 2000).

Next, we consider the overlap assumption, also known as positivity, in the current setting. We define the probability *density* of treatment realization w at time t given the history, $e_t(w) = f(W_t = w | \overline{H}_{t-1})$,

as the propensity score at time period t . Also, let f_h denote the probability density function of the stochastic intervention F_h . The assumption requires the ratio of propensity score over the density for the stochastic intervention, rather than the propensity score itself, is bounded away from zero.

Assumption 2 (Bounded relative overlap). *There exists a constant $\delta_W > 0$ such that $e_t(w) > \delta_W \cdot f_h(w)$ for all $w \in \mathcal{W}$.*

Assumption 2 ensures that all the treatment patterns which are possible under the stochastic intervention of interest can also be observed. This assumption enforces that the support of the intervention distribution has to be included in the support of the propensity score, and does not allow for interventions that assign positive mass to fixed treatments w .

4.2 The Propensity Score for Point Process Treatments

The propensity score plays an important role in our estimation. Here, we show that the propensity score for point process treatments has two properties analogous to those of the standard propensity score (Rosenbaum and Rubin, 1983). That is, the propensity score is a balancing score, and under Assumption 1 the treatment assignment is unconfounded conditional on the propensity score.

Proposition 1. *The propensity score $e_t(w)$ is a balancing score. That is, $f(W_t = w \mid e_t(w), \overline{H}_{t-1}) = f(W_t = w \mid e_t(w))$ holds for all t .*

In practice, Proposition 1 allows us to empirically assess the propensity score model specification by checking the predictive power of covariates in \overline{H}_{t-1} for the treatment W_t conditional on the propensity score. For example, if a covariate significantly improves prediction in a point process model of W_t after adjusting for the estimated propensity score, then the covariate is not balanced and the propensity score model is likely to be misspecified.

Proposition 2. *Under Assumption 1, the treatment assignment at time t is unconfounded given the propensity score at time t , that is, given $f(W_t \mid \overline{\mathbf{W}}_{t-1}, \overline{\mathcal{Y}}_T, \overline{\mathcal{X}}_T) = f(W_t \mid \overline{H}_{t-1})$, we have*

$$f(W_t \mid \overline{\mathbf{W}}_{t-1}, \overline{\mathcal{Y}}_T, \overline{\mathcal{X}}_T) = f(W_t \mid e_t(W_t)).$$

Proposition 2 shows that the potentially high-dimensional sets, \overline{H}_{t-1}^* and \overline{H}_{t-1} , can be reduced to the one dimensional propensity score $e_t(w)$ as a conditioning set sufficient for estimating the causal effects of W_t .

4.3 The Estimators

To estimate the causal estimands defined in Section 3, we propose propensity-score-based estimators that combine the inverse probability of treatment weighting (IPW) with the kernel smoothing of spatial point patterns. The estimation proceeds in two steps. First, at each time period t , the surface of outcome-active locations is spatially smoothed according to a chosen kernel. Then, this surface is weighted by the relative density of the observed treatment pattern under the stochastic intervention of interest and under the actual data generating process.

An alternative approach would be the direct modelling of the outcome. For example, one would model the outcome point process as a function of the past history following the g -computation in the standard longitudinal settings (Robins, 1986). However, such an approach would require an accurate specification of spatial spillover and temporal carryover effects. This is a difficult task in many applications. Instead, we focus on modelling the treatment assignment mechanism.

Formally, consider a univariate kernel $K : [0, \infty) \rightarrow [0, \infty)$ satisfying $\int K(u)du = 1$. Let K_b denote the scaled kernel defined as $K_b(u) = b^{-1}K(u/b)$ with bandwidth parameter b . We define $\widehat{Y}_t(F_h^M) : \Omega \rightarrow \mathbb{R}^+$ as

$$\widehat{Y}_t(F_h^M; \omega) = \prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)} \left[\sum_{s \in S_{Y_t}} K_b(\|\omega - s\|) \right], \quad (6)$$

where $\|\cdot\|$ denotes the Euclidean norm. The summation represents the spatially-smoothed version of the outcome point pattern at time period t . The product of ratios represents a weight similar to those in the marginal structural models (Robins *et al.*, 2000), but in accordance with the stochastic intervention F_h^M : each of the M terms represents the likelihood ratio of treatment W_j in the counterfactual world of the intervention F_h versus the actual world with the observed data at a specific time period.

Assuming that the kernel K is continuous, the estimator given in Equation (6) defines a continuous surface over Ω . The continuity of $\widehat{Y}(F_h^M)$ allows us to use it as an intensity function when estimating causal quantities. This leads to the following estimator for the expected number of outcome-active locations in any region B at time t , defined in Equation (2),

$$\widehat{N}_{Bt}(F_h^M) = \int_B \widehat{Y}_t(F_h^M; \omega) d\omega. \quad (7)$$

We can now construct the following estimator for the temporally-expected average potential outcome defined in Equation (4),

$$\widehat{N}_B(F_h^M) = \frac{1}{T - M + 1} \sum_{t=M}^T \widehat{N}_{Bt}(F_h^M). \quad (8)$$

We estimate the causal contrast between two interventions $F_{h_1}^M$ and $F_{h_2}^M$ defined in Equation (5) as,

$$\widehat{\tau}_B(F_{h_1}^M, F_{h_2}^M) = \widehat{N}_B(F_{h_2}^M) - \widehat{N}_B(F_{h_1}^M). \quad (9)$$

An alternative estimator of $N_{Bt}(F_h^M)$ could be obtained by replacing the kernel-smoothed version of the outcome in Equation (6) with the number of observed outcome active locations in B at time t . Even though this estimator has the same asymptotic properties discussed below, the kernel-smoothing of the outcome ensures that, for a specific intervention F_h^M , once the surface in Equation (6) is calculated, it can then be used to estimate the temporally-expected effects defined in Section 3 for any $B \subset \Omega$. In addition, it allows for the visualization of the outcome surface under an intervention, making it easier to identify the areas of increased or decreased activity as illustrated in Section 7.

In the next section we establish the asymptotic properties of the proposed IPW estimators. In our simulations (Section 6) and empirical study in (Section 7), we also use the Hájek estimator, which

standardizes the IPW weights and replaces Equation (8) with

$$\widehat{N}_B(F_h^M) = \sum_{t=M}^T \widehat{N}_{Bt}(F_h^M) / \sum_{t=M}^T \left\{ \prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)} \right\}. \quad (10)$$

We find that this Hájek estimators outperform the corresponding IPW estimators in finite samples, mirroring the existing results under other settings (e.g., Liu *et al.*, 2016; Cole *et al.*, 2021).

4.4 The Asymptotic Properties of the Proposed IPW Estimators

Below, we establish the asymptotic properties of the proposed IPW estimators. Our results differ from the existing asymptotic normality results in the causal inference literature in several ways. First, our inference is based on a single time series of point patterns that are both spatially and temporally dependent. Second, we employ a kernel-smoothed version of the outcome. Third, using martingale theory, we derive a new central limit theorem in time-dependent, observational settings. We now present the main theoretical results. All proofs are given in Appendix B.

Theorem 1 (Asymptotic Normality). *Suppose that Assumptions 1 and 2 as well as the regularity conditions (Assumption A.1) hold. Then, if the bandwidth $b_T \rightarrow 0$ and as $T \rightarrow \infty$, we have that*

$$\sqrt{T}(\widehat{N}_B(F_h^M) - N_B(F_h^M)) \xrightarrow{d} \mathcal{N}(0, v),$$

where v represents the probability limit of $(T - M + 1)^{-1} \sum_{t=M}^T v_t$ as $T \rightarrow \infty$ with

$$v_t = \text{Var} \left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)} N_B(Y_t) \mid \overline{H}_{t-M}^* \right] \quad \text{for } t \geq M.$$

The key idea of our proof is to separate the estimation error arising due to the treatment assignment W_t given the complete history \overline{H}_{t-1}^* , from the error due to spatial smoothing. Using martingale theory, we show that the former is \sqrt{T} -asymptotically normal, where the temporal dependence is controlled based on Assumption 1. The latter is shown to converge to zero at a rate faster than $1/\sqrt{T}$.

According to Theorem 1, the knowledge of v would enable asymptotic inference about the temporally-

expected potential outcome. The variance v is the converging point of $(T - M + 1)^{-1} \sum_{t=M}^T v_t$ where v_t represents a time period-specific variance. Unfortunately, since we only observe one treatment path for each time period t , we cannot directly estimate the time-specific variances, v_t , and thus v , without additional assumptions.

We circumvent this problem by using an upper bound of v , a quantity which we can consistently estimate. Specifically, let $v_t^* = \mathbb{E} \left\{ [\widehat{N}_{Bt}(F_h^M)]^2 \mid \overline{H}_{t-M}^* \right\}$. For v^* such that $(T - M + 1)^{-1} \sum_{t=M}^T v_t^* \xrightarrow{p} v^*$, we have $v \leq v^*$. Then, an α -level confidence interval for $N_B(F_h^M)$ based on the asymptotic variance bound v^*/T will achieve the nominal asymptotic coverage. Although v^* cannot be directly calculated either, there exists a consistent estimator of its upper bound, as stated in the following lemma:

Lemma 1 (Consistent Estimation of Variance Upper bound). *Suppose that Assumptions 1 and 2 and the regularity conditions (Assumption A.1) hold. Then, as $b_T \rightarrow 0$ and $T \rightarrow \infty$, we have*

$$\frac{1}{T - M + 1} \sum_{t=M}^T \left[\widehat{N}_{Bt}(F_h^M)^2 - v_t^* \right] \xrightarrow{p} 0.$$

In Appendix B.3 we extend the above results to the estimator $\widehat{\tau}_B(F_{h_1}^M, F_{h_2}^M)$.

So far, all of the theoretical results presented above have been established with the true propensity score $e_t(w)$. However, in practice, the propensity score is unknown and must be estimated. The next theorem shows that, when the propensity score is estimated under the correct model specification, the proposed estimator maintains its consistency and asymptotic normality. To prove this result, we extended classic M-estimation theory to multivariate martingale difference series, established a new central limit theorem for time series data, and derived the properties of the propensity score models under the spatio-temporal settings. To our knowledge, these results are new even though related results exist under the continuous time setting (Küchler *et al.*, 1999; Crimaldi and Pratelli, 2005). We believe that our results may be useful when studying the asymptotic properties of causal estimators in other dependent, observational settings (see Appendix B.4 for more details).

Theorem 2 (Asymptotic Normality Using the Estimated Propensity Score). *Suppose that Assumptions 1 and 2 as well as the regularity conditions (Assumptions A.1, A.2, A.3) hold. If the bandwidth $b_T \rightarrow 0$,*

then as $T \rightarrow \infty$, we have

$$\sqrt{T}(\widehat{N}_B(F_h^M) - N_B(F_h^M)) \xrightarrow{d} \mathcal{N}(0, v^e).$$

Next, we show that using the estimated propensity scores from a correctly specified model yields more efficient estimates than using the true propensity scores. This generalizes the well-known analogous result proved for the independent and identically distributed setting (e.g., [Hirano et al., 2003](#)) to the spatially and temporally dependent case (see [Zeng et al., 2021](#), for a similar result in a different dependent setting). Thus, even with the estimated propensity score, we can make asymptotically conservative inference based on the variance upper bound derived above.

Theorem 3 (Asymptotic Efficiency under the Estimated Propensity Score). *The estimator in Equation (8) based on the estimated propensity score from a correctly specified parametric model has asymptotic variance that is no larger than the asymptotic variance of the same estimator using the known propensity score. That is, for v in Theorem 1 and v^e in Theorem 2, we have $v^e \leq v$.*

The asymptotic results presented here require the area of interest, Ω , to be fixed while the number of time periods T increases. We note that point pattern treatments and outcomes might also arise in situations where the number of time periods T is fixed, but the area under study Ω grows to include more regions. In [Appendix B.5](#), we provide an alternative causal inference framework for point pattern treatments under this new design by extending our causal estimands, estimation and asymptotic results to the spatio-temporal setting with an area consisting of a growing number of independent regions.

5 Sensitivity Analysis

The validity of our estimators critically relies upon the assumption of no unmeasured confounding (Assumption 1). We develop a sensitivity analysis to address the potential violation of this key identification assumption. Specifically, we extend the sensitivity analysis pioneered by [Rosenbaum \(2002\)](#) to the spatio-temporal context and to the Hájek estimator with standardized weights, which we consider in our simulation and empirical studies.

Suppose there exists an unmeasured, potentially time-varying confounder U_t . We assume that the unconfoundedness assumption holds only after conditioning on the realized history of this unobserved confounder as well as \overline{H}_{t-1} , i.e.,

$$f(W_t | \overline{\mathbf{W}}_{t-1}, \overline{\mathcal{Y}}_T, \overline{\mathcal{X}}_T, \overline{\mathcal{U}}_T) = f(W_t | \overline{H}_{t-1}, \overline{U}_t),$$

where $\overline{\mathcal{U}}_T$ represents the collection of all potential values of U_t across all time points $t = 1, 2, \dots, T$ whereas \overline{U}_t represents the history of realized but unmeasured confounder U up to time t . Note that U_t can be correlated with the observed confounders.

The existence of an unmeasured confounder invalidates the inference based on the propensity score with observed covariates alone $e_t(w)$ because the true propensity score, denoted by $e_t^*(w) = f_{W_t}(w | \overline{H}_{t-1}, \overline{U}_t)$, conditions on the history of the unmeasured confounder \overline{U}_t . To develop a sensitivity analysis, we assume the ratio of estimated versus true propensity scores for the realized treatment W_t is bounded by a value $\Gamma (\geq 1)$,

$$\frac{1}{\Gamma} \leq \rho_t = \frac{e_t(W_t)}{e_t^*(W_t)} \leq \Gamma.$$

A larger value of Γ allows a greater degree of violation of the unconfoundedness assumption.

In our application, we use the Hájek-version of the proposed estimator, which we find to be more stable than the IPW estimator (see Section 6). Thus, to develop a sensitivity analysis, we derive an algorithm for bounding the Hájek estimator for stochastic interventions for each fixed value of Γ (see Appendix D for the sensitivity analysis of the IPW estimator). Specifically, for all values of $\boldsymbol{\rho} = (\rho_1, \rho_2, \dots, \rho_T) \in [\Gamma^{-1}, \Gamma]^T$, we wish to bound the following two quantities:

$$\begin{aligned} \widehat{N}_{\boldsymbol{\rho}}(F_h) &= \frac{\sum_{t=1}^T \rho_t w_t(F_h) \widetilde{N}_B(Y_t)}{\sum_{t=1}^T \rho_t w_t(F_h)}, \quad \text{and} \\ \widehat{\tau}_{\boldsymbol{\rho}}(F_{h_1}, F_{h_2}) &= \frac{\sum_{t=1}^T \rho_t w_t(F_{h_2}) \widetilde{N}_B(Y_t)}{\sum_{t=1}^T \rho_t w_t(F_{h_2})} - \frac{\sum_{t=1}^T \rho_t w_t(F_{h_1}) \widetilde{N}_B(Y_t)}{\sum_{t=1}^T \rho_t w_t(F_{h_1})} \end{aligned}$$

where

$$w_t(F_h) = \frac{f_h(W_t)}{e_t(W_t)} \quad \text{and} \quad \tilde{N}_B(Y_t) = \int_B \sum_{s \in S_{Y_t}} K_b(\|\omega - s\|) d\omega.$$

Below we show how to formulate the bounding problem for $\hat{N}_\rho(F_h)$ as a linear program, and how to use the bounds for $\hat{N}_\rho(F_h)$ to also bound the effect estimator $\hat{\tau}_\rho(F_{h_1}, F_{h_2})$.

Theorem 4 (Bounding the Causal Quantities).

1. The problem of maximizing $\hat{N}_\rho(F_h)$ over $\rho \in [\Gamma^{-1}, \Gamma]^T$ is equivalent to the following linear program,

$$\begin{aligned} & \text{maximize}_{\rho^*} \sum_{t=1}^T \rho_t^* w_t(F_h) \tilde{N}_B(Y_t) \\ & \text{subject to} \quad \frac{\kappa}{\Gamma} \leq \rho_t^* \leq \Gamma \kappa, \quad \sum \rho_t^* w_t(F_h) = 1, \quad \text{and} \quad \kappa \geq 0. \end{aligned}$$

where $\rho = \rho^* / \kappa$.

2. Suppose that ρ_j^{\max} and ρ_j^{\min} represent the values of ρ that maximize and minimize $\hat{N}_\rho(F_{h_j})$, respectively, for $j = 1, 2$. Then, the bounds for the causal effect are obtained as,

$$\hat{N}_{\rho_2^{\min}}(F_{h_2}) - \hat{N}_{\rho_1^{\max}}(F_{h_1}) \leq \hat{\tau}_\rho(F_{h_1}, F_{h_2}) \leq \hat{N}_{\rho_2^{\max}}(F_{h_2}) - \hat{N}_{\rho_1^{\min}}(F_{h_1}). \quad (11)$$

The proof for bounding $\hat{N}_\rho(F_h)$ is based on the Charnes-Cooper transformation of linear fractionals (Charnes and Cooper, 1962), and the proof for bounding $\hat{\tau}_\rho(F_{h_1}, F_{h_2})$ is given in Appendix D. For bounding $\hat{N}_\rho(F_h)$, this proposition allows us to use a standard linear algorithm to obtain the optimal solution for (ρ^*, κ) and transform it back to the optimal solution ρ . Then, we can use these bounds to also acquire bounds on the effect estimator. Since all bounds are wider for a greater value of Γ , the estimated effects are robust to propensity score misspecification up to the smallest value of Γ for which the interval of bounds in Equation (11) includes 0. Due to the standardization of weights in the Hájek estimator, the bound in Equation (11) is conservative, in the sense that, if the causal estimate is shown to be robust up to some value Γ , then it is robust up to an even greater degree of propensity score model misspecification $\Gamma^* \geq \Gamma$. Similar bounds can be derived for the stochastic interventions that take place over multiple time periods (see Appendix D for details).

The propensity score modelling in our spatio-temporal setting is much more complex with an infinite number of potential treatment locations than in the conventional cross-sectional setting. As a result, the modelling uncertainty for the propensity score is much greater. This renders the scale of Γ in the two settings not directly comparable, and the value of Γ is expected to be much closer to the null value of one in the spatio-temporal context.

6 Simulation Studies

We conduct simulation studies to empirically investigate several key theoretical properties of the proposed methodology: (a) the performance of our estimator under different stochastic interventions and as the number of time periods increases, (b) the accuracy of the asymptotic approximation, (c) the difference between the theoretical variance bound and the actual variance, (d) the performance of the inferential approach based on the estimated asymptotic variance bound, (e) the relative efficiency of the estimator when using the true and estimated propensity scores, and (f) the balancing properties of the estimated propensity score. We use the `spatstat` R package (Baddeley *et al.*, 2015) to generate point patterns from Poisson processes and fit Poisson process models to the simulated data.

6.1 The Study Design

To construct a realistic simulation design, we base our data generating process on the observed data from our application. We consider a time series of point patterns of length $T \in \{200, 400, 500\}$. For each time series length T , 200 data sets are generated. The scenario with $T = 500$ closely resembles our observed data, which have $T = 469$.

Time-varying and time-invariant confounders. Our simulation study includes two time-invariant and two time-varying confounders. We base the first time-invariant confounder on the distance from Iraq’s road network and its borders, by defining its value at location $\omega \in \Omega$ as $X^1(\omega) = \exp\{-3D_1(\omega)\} + \log(D_2(\omega))$, where $D_1(\omega)$ is the distance from ω to the closest road, and $D_2(\omega)$ is the distance to the country’s border. This covariate is shown in Figure 3a. The second covariate is defined similarly as



(a) Time-invariant confounder $X^1(\omega)$ (b) Realization of time-varying confounder $X_t^3(\omega)$ (c) Distribution of treatment point patterns

Figure 3: Simulated Data. Panel (a) shows one of the two time-invariant confounders representing the exponential decay of distance to the road network. Panel (b) shows one realization for one of the time-varying confounders. After generating points from a non-homogeneous Poisson process, which depends on the observed airstrike density, we define the time-varying confounder as the exponential decay of distance to these points. Panel (c) shows the estimated log-density of treatment patterns, which is used as the density ϕ in the definition of stochastic interventions.

$X^2(\omega) = \exp\{-D_3(\omega)\}$ where $D_3(\omega)$ is the distance from ω to Baghdad.

We generate the time-varying confounders, $X_t^3(\omega)$ and $X_t^4(\omega)$, using the kernel-smoothed density of the observed airstrike and attack patterns. Specifically, we pool all airstrike locations across time and estimate the density of airstrike patterns, denoted by $\hat{f}(\omega)$ at location ω (shown in the right plot of Figure 6). Based on this density, we draw a point pattern from a non-homogeneous Poisson point process with intensity function $\lambda^{X^3}(\omega) = \exp\{\rho_0 + \rho_1 \hat{f}(\omega)\}$, for $\rho_0 \approx -2.7$ and $\rho_1 = 8$, and define $X^3(\omega)$ as $\exp\{-D_4(\omega)\}$, where $D_4(\omega)$ is the distance from location ω to the closest point. We generate $X^4(\omega)$ similarly based on the estimated density for insurgent attacks, and for corresponding values $\rho_0 \approx -3.2$ and $\rho_1 = 7$. Figure 3b shows one realization of $X_t^3(\omega)$.

Spatio-temporal point processes for treatment and outcome variables. For each time period $t \in \mathcal{T}$, we generate W_t from a non-homogeneous Poisson process that depends on all confounders $\mathbf{X}_t(\omega) = (X^1(\omega), X^2(\omega), X_t^3(\omega), X_t^4(\omega))^\top$, as well as the previous treatment and outcome realizations, W_{t-1} and Y_{t-1} . The intensity of this process is given by

$$\lambda_t^W(\omega) = \exp\{\alpha_0 + \boldsymbol{\alpha}_X^\top \mathbf{X}_t(\omega) + \alpha_W W_{t-1}^*(\omega) + \alpha_Y Y_{t-1}^*(\omega)\}, \quad (12)$$

where $W_{t-1}^*(\omega) = \exp\{-2D_W(\omega)\}$ and $Y_{t-1}^*(\omega) = \exp\{-2D_Y(\omega)\}$ with $D_W(\omega)$ and $D_Y(\omega)$ being the minimum distance from ω to the points in $S_{W_{t-1}}$ and $S_{Y_{t-1}}$, respectively.

Similarly, we generate Y_t from a non-homogeneous Poisson process with intensity

$$\lambda_t^Y(\omega) = \exp\{\gamma_0 + \gamma_X^\top \mathbf{X}_t(\omega) + \gamma_2 X_{t-1}^2(\omega) + \gamma_W W_{(t-3):t}^*(\omega) + \gamma_Y Y_{t-1}^*(\omega)\}, \quad (13)$$

where $W_{(t-3):t}^*(\omega) = \exp\{-2D_W^*(\omega)\}$ with $D_W^*(\omega)$ being the distance from ω to the closest points in $\bigcup_{j=t-3}^t S_{W_j}$. This specification imposes a lag-three dependence of the outcome on the lagged treatment process. The model leads to an average of 5.5 treatment-active locations and 31 outcome-active locations within each time period, resembling the frequency of events in our observed data. The spatial distribution of generated treatment point patterns also resembles the observed one (compare Figure 3c to the right plot of Figure 6). The simulated and observed outcome point patterns also have similar distributions.

Stochastic interventions. We consider stochastic interventions of the form F_h^M for a non-homogeneous Poisson process with intensity h , which is defined as $h(\omega) = c\phi(\omega)$ for c ranging from 3 to 8, and surface ϕ set to the density shown (in logarithm) in Figure 3c. This definition of stochastic intervention based on the treatment density aligns with the specification in our study in Section 7. We consider varying the intervention duration by setting $M \in \{1, 3, 7, 30\}$. We also examine lagged interventions over three time periods, i.e., $F_h = F_{h_3} \times F_{h_2} \times F_{h_1}$. The intervention for the first time period F_{h_3} is a Poisson process with intensity $h_3(\omega) = c\phi(\omega)$ for c ranging from 3 to 7, whereas $F_{h_2} = F_{h_1}$ is a non-homogeneous Poisson process with intensity $5\phi(\omega)$. For each stochastic intervention, we consider three regions of interest, B , of different sizes, representing the whole country, the Baghdad administrative unit, and a small area in northern Iraq which includes the town of Mosul.

Approximating the true estimands. Equation (13) shows that the potential outcomes depend on the realized treatments during the last four time periods as well as the realized outcomes from the previous time period. This implies that the estimands for all interventions, even for $M > 4$, depend on the ob-

served treatment and outcome paths and are therefore not constant across simulated data sets. Therefore, we approximate the true values of the estimands in each data set in the following manner. For each time period t , and each $r = 1, 2, \dots, R$ repetition, we generate realizations $w_{t-M+1}^{(r)}, \dots, w_{t-1}^{(r)}, w_t^{(r)}$ from the intervention distribution $F_{\mathbf{h}}$. Based on the treatment path $(\overline{\mathbf{W}}_{t-M}, w_{t-M+1}^{(r)}, \dots, w_t^{(r)})$, we generate outcomes $y_{t-M+1}^{(r)}, \dots, y_t^{(r)}$ using Equation (13). This yields $S_{y_t}^{(r)}$, which contains the outcome-active locations based on one realization from the stochastic intervention. Repeating this process R times and calculating the average number of points that lie within B provides a Monte Carlo approximation of $N_{Bt}(F_{\mathbf{h}})$, and further averaging these over time gives an approximation of $N_B(F_{\mathbf{h}})$.

Estimation. We estimate the expected number of points $N_B(F_{\mathbf{h}})$ and the effect of a change in the intervention on this quantity $\tau_B(F_{\mathbf{h}'}, F_{\mathbf{h}''})$ using the following estimators: (a) the proposed estimators defined in Equations (8) and (9) with the true propensity scores; (b) the same proposed estimators with the estimated propensity scores based on the correctly-specified model; (c) the above two estimators with the Hájek-type standardization in (10); and (d) the unadjusted estimator based on the propensity score model using a homogeneous Poisson process with no predictor.

All estimators utilize the smoothed outcome point pattern. Spatial smoothing is performed using Gaussian kernels with standard deviation equal to $10T^{-2/3}\delta$, which is decreasing in T , and for δ scaling the bandwidth according to the size of the geometry under study. We choose this bandwidth such that for $T = 500$ (the longest time series in our simulation scenario) the bandwidth is approximately 0.5, slightly smaller than the size of the smallest region of interest B (square with edge equal to 0.75). We discuss the choice of the bandwidth in Section 7.3.

Theoretical variance and its upper bound. Theorem 1 provides the expression for the asymptotic variance of the proposed IPW estimator. We compute Monte Carlo approximations to this variance and its upper bound. Specifically, for each time period t and each replication r , the computation proceeds as follows: 1) we generate treatment and outcome paths $w_{t-M+1}^{(r)}, y_{t-M+1}^{(r)}, \dots, w_t^{(r)}, y_t^{(r)}$ using the distributions specified in Equations (12) and (13), 2) using the data $(w_{t-M+1}^{(r)}, \dots, w_t^{(r)})$ and the outcome $y_t^{(r)}$,

we compute the estimator according to Equations (6) and (7), and finally 3) we calculate the variance and the second moment of these estimates over R replications, which can be used to compute the asymptotic variance and variance bound of interest. Their averages over time give the desired Monte Carlo approximations. We use a similar procedure to approximate the theoretical variance and variance bound of $\widehat{\tau}_B(F_{\mathbf{h}'}, F_{\mathbf{h}''})$.

Estimating the variance bound and the resulting inference. We use Lemma 1 to estimate the variance bound. This estimated variance bound is then used to compute the confidence intervals and conduct a statistical test of whether the causal effect is zero. Inference based on the Hájek estimator is discussed in Appendix C.

Balancing property of the propensity score. Using the correctly specified model, we estimate the propensity score at each time period t . The inverse of the estimated propensity score is then used as the weight in the weighted Poisson process model for W_t with the intensity specified in Equation (12). We compare the statistical significance of the predictors between the weighted and unweighted model fits. Large p -values under the weighted model would suggest that the propensity score adequately balances the confounding variables.

Relative efficiency of estimators based on the true and estimated propensity score. According to Theorem 3, the asymptotic variance of the estimator based on the true propensity score is at least as large as that of the estimator based on the estimated propensity score. We investigate the relative magnitude of the Monte Carlo approximations of the corresponding two variances.

6.2 Simulation Results

Figure 4 presents the results for all the stochastic interventions that were considered. The top four rows show how the (true and estimated) average potential outcomes in the whole region ($B = \Omega$) change as the intensity varies under interventions F_h^M for $M \in \{1, 3, 7, 30\}$, respectively. The last row shows how the true and estimated average potential outcomes in the same region change under the three time period

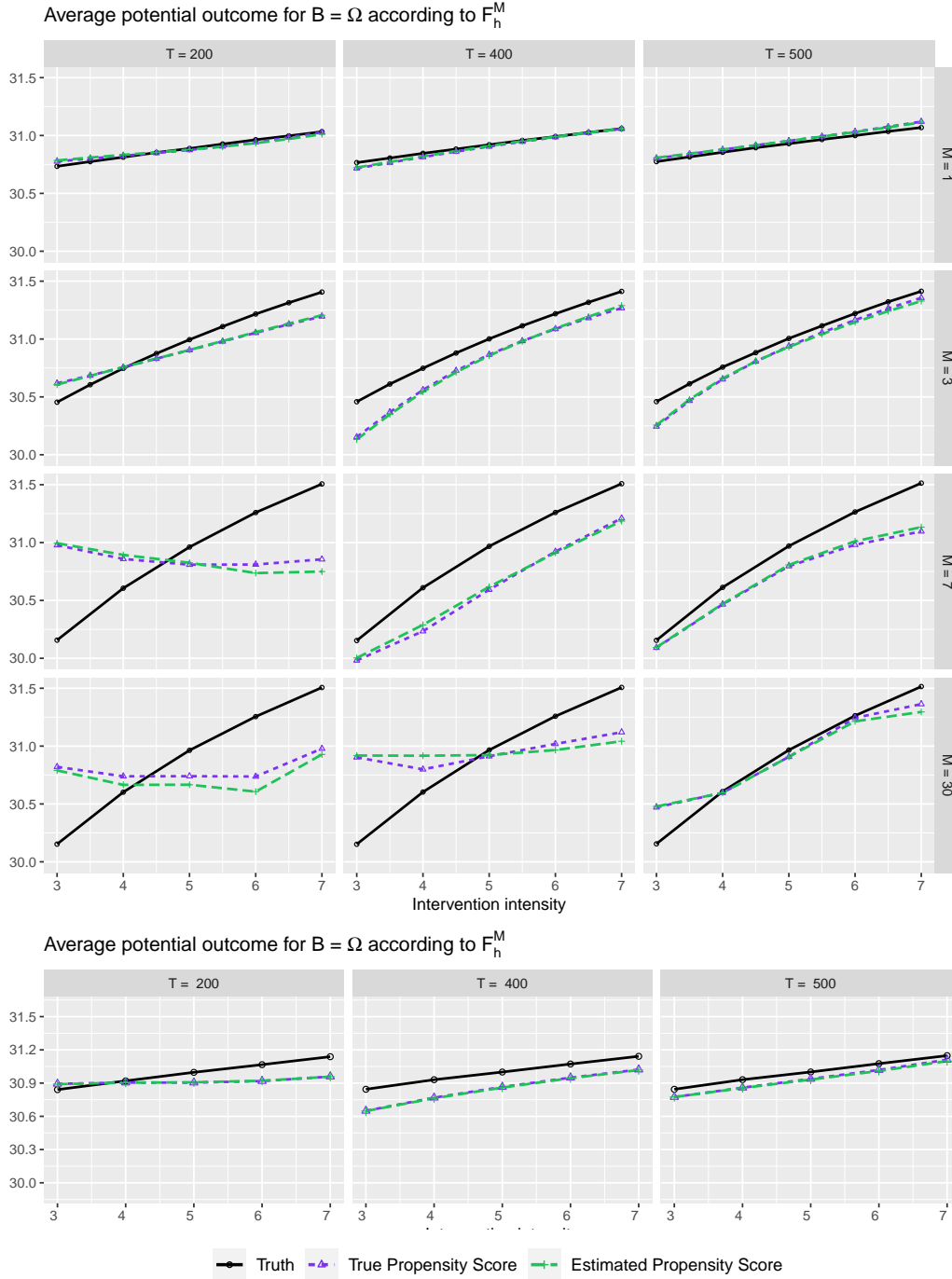


Figure 4: Simulation Results for the Average Potential Outcomes. In the top four rows, we present the true and estimated average potential outcomes in $B = \Omega$ under interventions F_h^M with the varying intensity (horizontal axis) and $M \in \{1, 3, 7, 30\}$ (rows), respectively. In the bottom row, we consider the average potential outcome for the lagged intervention over three time periods F_h , with the varying intensity of F_h shown on the horizontal axis. The black lines with solid circles represent the truths, the Hájek estimator based on the true propensity score is shown in purple, and the Hájek estimator based on the estimated propensity score is in green.

lagged interventions when the intensity at three time periods ago ranges from 3 to 7. For both simulation scenarios, we vary the length of the time series from 200 (left column) to 500 (right column).

The unadjusted estimator returned values that are too far from the truth and are not shown here. We find that the accuracy of the proposed estimator improves as the number of time periods increases. Notice that the convergence is slower for larger values of M . This is expected because the uncertainty of the treatment assignment is greater for a stochastic intervention with a longer time period. We find that the Hájek estimator performs well across most simulation scenarios even when T is relatively small and M is large. The IPW estimator (investigated more thoroughly in Appendix F) tends to suffer from extreme weights because the weights are multiplied over the intervention time periods as shown in Equation (6). These results indicate a deteriorating performance of the IPW estimator as the value of M increases, whereas the standardization of weights used in the Hájek estimator appears to partially alleviate this issue. Results were comparable for the two other sets B .

Next, we compare the true theoretical variance, v/T , with the variance bound v^*/T and its consistent estimator (see Lemma 1). We assess the conservativeness of the theoretical variance bound by focusing on the proposed estimators with the true propensity score. Figure 5 shows the results of an intervention F_h^M for $M \in \{1, 3\}$, and for region $B = \Omega$. The results for the other regions are similar and hence omitted.

First, we focus on the theoretical variance and variance bound (blue line with open circles, and orange dotted lines with open triangles, respectively). As expected, the true variance decreases as the total number of time periods increases, and the theoretical variance bound is at least as large as the variance. In the setting with $M = 3$, the theoretical variance follows the variance closely, evident by the fact that the two lines are essentially indistinguishable. We have found this to be the case in all scenarios with higher uncertainty, indicating that the theoretical variance bound is not overly conservative. Indeed, the variance bound is visibly larger than the true variance only in the low-variance scenarios of interventions over a single time period, as shown in the top row of Figure 5 (and in Appendix E.1).

Standard deviation and bound of the average potential outcome estimator for $B = \Omega$

Intervention over 1 time point

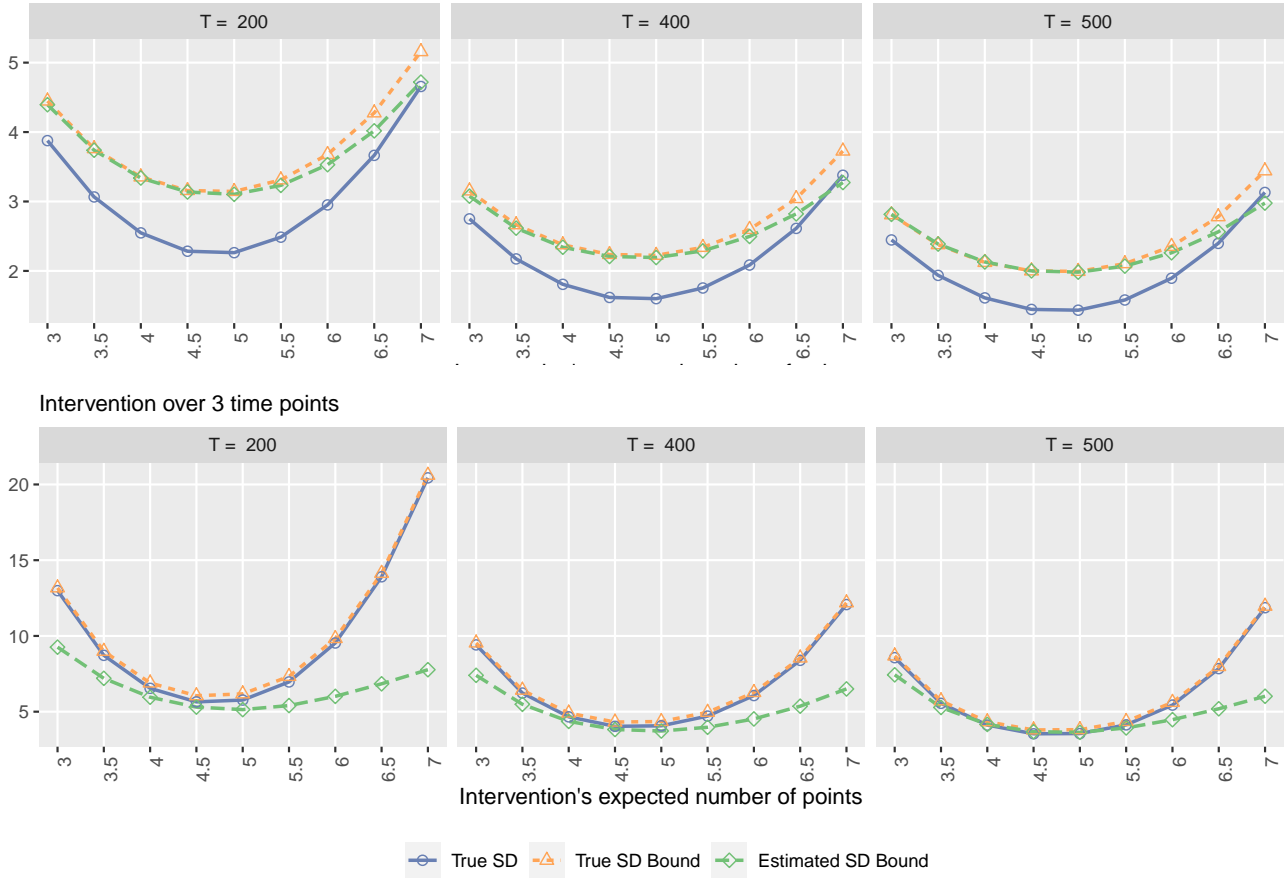


Figure 5: Comparison of the Asymptotic Standard Deviation with the True and Estimated Asymptotic Standard Deviation Bound for the Average Potential Outcome Estimator. The comparison is based on the varying number of expected points (horizontal axis) under the stochastic intervention F_h^M taking place over $M = 1$ (top row) and $M = 3$ (bottom row) time periods. The columns correspond to a simulation setting with a different time series length.

Second, we compare the theoretical variance bound with the estimated variance bound (green dashed lines with open rhombuses). As the length of time series increases, the estimated variance bound more closely approximates its theoretical value (consistent with Lemma 1). Furthermore, the estimated variance bound is close to its theoretical value under low uncertainty scenarios and when the intervention intensity more closely resembles that of the actual data generating process. However, we find that the estimated variance bound underestimates the true variance bound in high uncertainty scenarios, and convergence to its true value is slower for larger values of M (see Appendix E.1).

We also compare the variance of the estimator based on the true propensity score with that of the estimator based on the estimated propensity score. Table 1 shows the ratio of the Monte Carlo variances which, according to Theorem 3, should be larger than 1, asymptotically. Consistent with the above simulation results, we find that the ratio is above 1 for interventions over one and three time periods. In addition, the ratio is largest in the low uncertainty scenarios where either the number of intervention periods, M , is small, or the expected number of points is near the observed value under the intervention ($c \approx 5$). In contrast, in the high uncertainty situations with longer intervention periods, e.g. $M \in \{7, 30\}$, the ratio remains below 1, implying that the asymptotic approximation may not be sufficiently accurate for the sample sizes considered.

In Appendices E.2 and E.3, we also investigate the performance of the inferential procedure based on the true variance, true variance bound, and estimated variance bound, for both the IPW and Hájek estimators. The confidence interval for the IPW estimator tends to yield coverage close to its nominal level only for the interventions over a small number of time periods. In contrast, the confidence interval for the Hájek estimator has good coverage probability even for the interventions over a larger number of time periods. Partly based on these findings, we use the Hájek estimator and its associated confidence interval in our empirical application (see Section 7).

Table 1: Variance Ratio of the Proposed Estimator based on the True Propensity Score over the Proposed Estimator based on the Estimated Propensity Score. The results are based on Monte Carlo approximation with $T = 500$. The estimated propensity score is obtained from the correctly specified model. If the ratio is greater than 1, the estimated propensity score yields more efficient estimator than the true propensity score. We consider interventions that are constant over all intervention time periods, F_h^M for $M \in \{1, 3, 7, 30\}$ (top four rows), and the lagged intervention over three time periods $F_h = F_{h_3} \times F_{h_2} \times F_{h_1}$ (bottom row).

	Expected number of treatment active locations under the intervention				
	$c = 3$	$c = 4$	$c = 5$	$c = 6$	$c = 7$
$M = 1$	1.24	1.38	1.41	1.32	1.08
$M = 3$	1.09	1.18	1.24	1.14	1.06
$M = 7$	1.07	0.85	1.08	0.61	0.54
$M = 30$	0.60	0.75	0.87	0.58	0.75
Lagged	1.08	1.21	1.24	1.20	1.11

Finally, we evaluate the performance of the propensity score as a balancing score (Proposition 1). In Appendix E.4, we show that the p-values of the previous outcome-active locations variable (Y_{t-1}^* in Equation (12)) are substantially greater in the weighted propensity score model than in the unweighted model, where the weights are equal to the inverse of the estimated propensity score. These findings are consistent with the balancing property of the propensity score.

In Appendix F we present an alternative simulation study, though all qualitative conclusions remain unchanged.

7 Empirical Analyses

In this section, we present our empirical analyses of the datasets introduced in Section 2. We first describe the airstrike strategies of interest and then discuss the causal effect estimates obtained under those strategies.

7.1 Airstrike Strategies and Causal Effects of Interest

We consider hypothesized stochastic interventions that generate airstrike locations based on a simple non-homogeneous Poisson point process with finite and non-atomic intensity $h : \Omega \rightarrow [0, \infty)$. We first specify a baseline probability density ϕ_0 over Ω . To make this baseline density realistic and increase the credibility of the overlap assumption, we use the airstrike data during January 1 – September 24, 2006 to define the baseline distribution ϕ_0 for our stochastic interventions. This subset of the data is not used in the subsequent analysis. The left plot of Figure 6 shows the estimated baseline density, using kernel-smoothing of airstrikes with an anisotropic Gaussian kernel and bandwidth specified according to Scott’s rule of thumb (Scott, 1992).

We consider the following three questions: (1) How does an increase in the number of airstrikes affect insurgent violence? (2) How does the shift in the prioritization of certain locations for airstrikes change the spatial pattern of insurgent attacks? (3) How long does it take for the effects of change in these airstrike strategies to be realized? The last question examines how quickly the insurgents respond

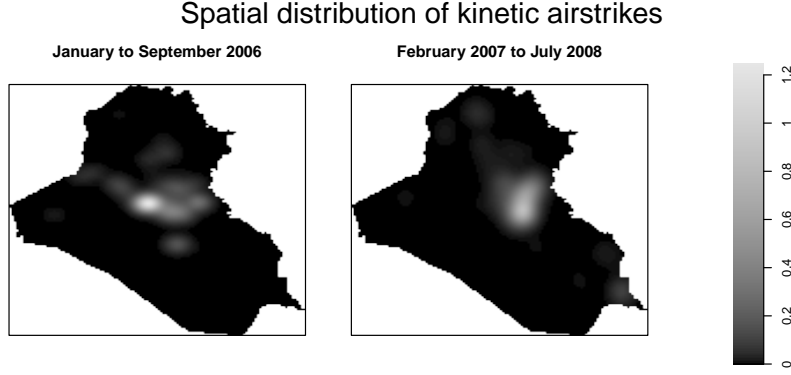


Figure 6: Spatial Density Estimate of Airstrike Locations during January 1 – September 24, 2006 (left) and the Entire Study Period February 2007 – July 2008 (right).

to the change in airstrike strategy.

We address the first question by considering stochastic interventions that have the same spatial distribution but vary in the expected number of airstrikes. We represent such strategies using intensities $h(\omega) = c\phi_0(\omega)$ with different values of $c > 0$. Since $\int_{\Omega} h(\omega)d\omega$ represents the expected number of points from a Poisson point process, these interventions have the same spatial distribution ϕ_0 , but the number of airstrikes monotonically increases as a function of c . In our analysis, we consider $\{1, 2, \dots, 6\}$ as the range of c which corresponds to the expected number of airstrikes per day, in agreement with the observed data.

For the second question, we fix the expected number of airstrikes but vary their focal locations. To do this, we specify a distribution over Ω with power-density $d_{\alpha}(\omega) = d(\omega)^{\alpha} / (\int_{\Omega} d(\omega)^{\alpha})$ and modes located at $s_f \in \Omega$. Based on d_{α} , we specify $h_{\alpha}(\omega) = c_{\alpha}\phi_0(\omega)d_{\alpha}(\omega)$ where c_{α} satisfies the constraint $\int_{\Omega} h_{\alpha}(\omega)d\omega = c$, so that the overall expected number of airstrikes remains constant. Locations in s_f are increasingly prioritized under h_{α} for increasing α . For our analysis, we choose the center of Baghdad to be the focal point s_f and d_{α} to be the normal distribution centered at s_f with precision α . We set the expected number of airstrikes per day c to be 3, and vary the precision parameter α from 0 to 3. The visualization of the spatial distributions in h_{α} for the different values of α is shown in Figure A.15.

As discussed in Section 3, for both of these questions, we can specify airstrike strategies of interest

taking place over a number of time periods, M , by specifying the stochastic interventions as $F_{\mathbf{h}} = F_h^M$. In addition, we may also be interested in the lagged effects of airstrike strategies as in the third question. We specify lagged intervention to be the one which differs only for the M time periods ago, i.e., $F_{\mathbf{h}} = F_{h_0}^{M-1} \times F_{h_1}$, where $h_0 = \phi_0$ represents the baseline intensity (with $c = 1$), and $h_1 = c\phi_0$ is the increased intensity with different values of c ranging from 1 to 6. We assume that insurgent attacks at day t do not affect airstrikes on the same day, and airstrikes at day t can only affect attacks during subsequent time periods. Thus, causal quantities for interventions taking place over M time periods refer to insurgent attacks occurring M days later. For our analysis, we consider values of M which correspond to 1 day, 3 days, 1 week, and 1 month.

Although full investigation is beyond the scope of this paper, in Appendix G.3, we briefly consider an extension to adaptive interventions over a single time period ($M = 1$), and discuss challenges when considering adaptive interventions over multiple time periods ($M > 1$).

7.2 The Specification and Diagnostics of the Propensity Score Model

Our propensity score model is a non-homogeneous Poisson point process model with intensity $\lambda_t(\omega) = \exp\{\beta^\top \mathbf{X}_t(\omega)\}$ where \mathbf{X} includes an intercept, temporal splines, and 32 spatial surfaces including all the covariates. The two main drivers of military decisions over airstrikes are the prior number and locations of observed insurgent attacks and airstrikes, which are expected to approximately satisfy unconfoundedness of Assumption 1. Our model includes the observed airstrikes and insurgent attacks during the last day, week, and month (6 spatial surfaces). For example, the airstrike history of time t during the previous week is $\overline{W}_{t-1}^*(\omega) = \sum_{j=1}^7 \sum_{s \in S_{W_{t-j}}} \exp\{-\text{dist}(s, \omega)\}$, which represents a surface on Ω with locations closer to the airstrikes in the previous week having greater values than more distant locations.

Our propensity score model also includes additional important covariates that might affect both airstrikes and insurgent attacks. We adjust for shows-of-force (i.e., simulated bombing raids designed to deter insurgents) that occurred one day, one week, and one month before each airstrike (3 spatial surfaces). Patterns of U.S. aid spending might also affect the location and number of insurgent attacks and

airstrikes, as we discussed in Section 2. We therefore include the amount of aid spent (in U.S. dollars) in each Iraqi district in the past month as a time-varying covariate (1 spatial surface). Finally, we also incorporate several time-invariant spatial covariates, including the airstrike’s distance from major cities, road networks, rivers, and the population (logged, measured in 2003) of the governorate in which the airstrike took place (4 spatial surfaces). Lastly, we include separate predictors for distances from local settlements in each of the Iraqi districts to incorporate any area specific effects (18 spatial surfaces).

We evaluate the covariate balance by comparing the p -values of estimated coefficients in the propensity score model to the p -values in the weighted version of the same model, where each time period is inversely weighted by its truncated propensity score estimate (truncated above at the 90th quantile). Although 13 out of 35 estimated coefficients had p -values smaller than 0.05 in the fitted propensity score model, all the p -values in the weighted propensity score model are close to 1, suggesting that the estimated propensity score adequately balances these confounders (see Figure A.14 of Appendix E.4).

7.3 The Choice of the Bandwidth Parameter for the Spatial Kernel Smoother

The kernel smoothing part of our estimator is not necessary for estimating the number of points within any set $B \subset \Omega$ since we can simply use an IPW estimator based on the observed number of points within B . However, kernel smoothing is useful for visualizing the estimated intensities of insurgent attacks under an intervention of interest over the entire country. One can also use it to acquire estimates of the expected number of insurgent attacks under the intervention for any region of Iraq by considering the intensity’s integral over the region. Theorem 1 shows that, for any set $B \subset \Omega$, kernel smoothing does not affect the estimator’s asymptotic normality as long as the bandwidth converges to zero. In practice, the choice of the bandwidth should be partly driven by the size of the sets B .

In our analysis, we estimate the causal quantities for the entire country and the Baghdad administrative unit. We choose an adaptive bandwidth separately for each outcome using the `spatstat` package in R. We consider all observed outcome event locations during our study period, and use Scott’s criterion for choosing an optimal, constant bandwidth parameter for isotropic kernel estimation (Scott, 1992).

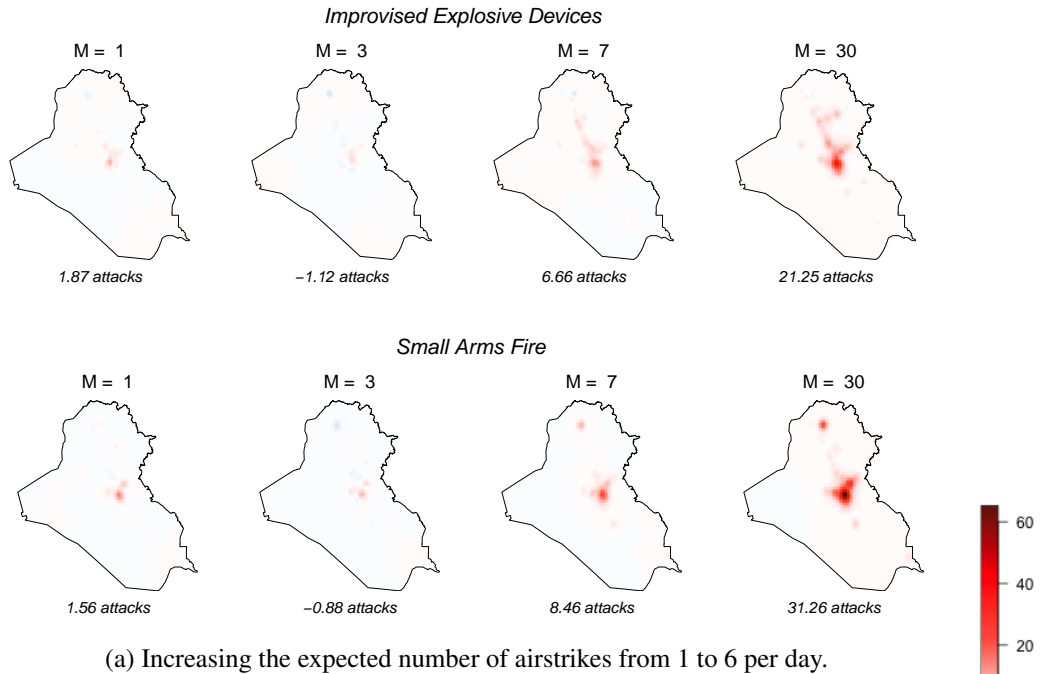
Using the estimated density as the pilot density, we calculate the optimal adaptive bandwidth surface according to Abramson’s inverse-square-root rule (Abramson, 1982). This procedure yields a value of the bandwidth used for kernel smoothing at each outcome event location.

7.4 Findings

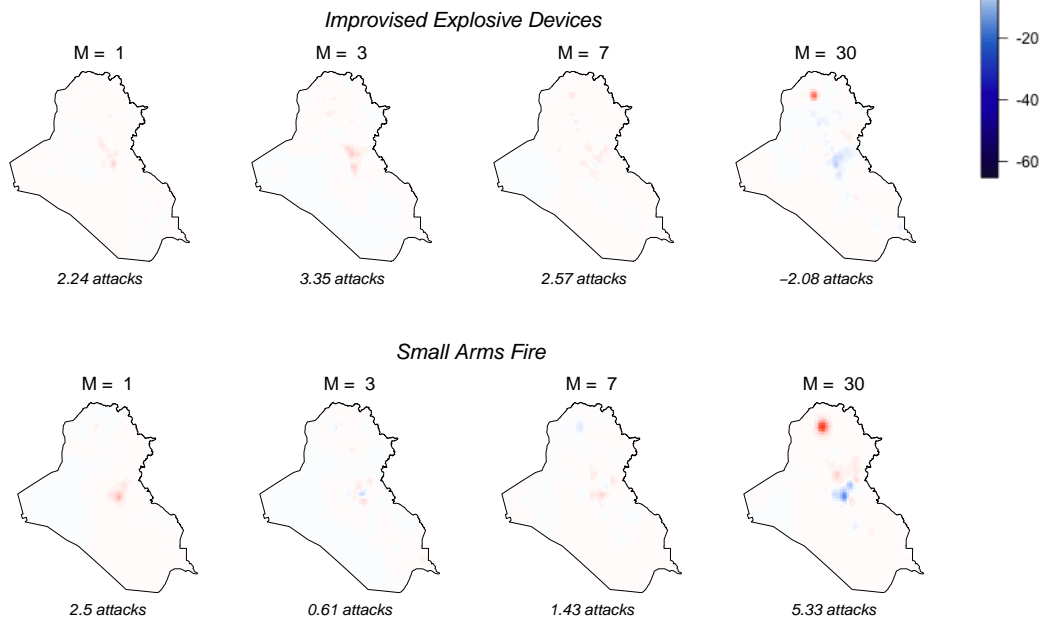
Figure 7 illustrates changes in the estimated intensity surfaces for insurgent attacks (measured using IEDs and SAFs) when increasing the expected number of airstrikes (the first two rows) and when shifting the focal point of airstrikes to Baghdad (the bottom two rows), with the varying duration of interventions, $M = 1, 3, 7, 30$ days (columns). These surfaces can be used to estimate the causal effect of a change in the intervention over any region. Dark blue areas represent areas where the change in the military strategy would reduce insurgent attacks, whereas red areas correspond to those with an increase in insurgent attacks. Statistical significance of these results is shown in Appendix G.2.

The figure reveals a number of findings. First, there is no substantial change in insurgent attacks if these interventions last only for one or three days. When increasing airstrikes for a longer duration, however, a greater number of insurgent attacks are expected to occur. These changes are concentrated in the Baghdad area and the roads that connect Baghdad and the northern city of Mosul. These patterns apply to both IEDs and SAFs with slightly greater effects estimated for SAFs. Under Assumption 1, these results indicate that, far from suppressing insurgent attacks, airstrikes actually *increase* them over time. In this setting, airstrikes are counterproductive, failing to reduce insurgent violence while also victimizing civilians.

The effect estimates for shifting the focal point of airstrikes to Baghdad for 1, 3, or 7 days are close to null. However, when the intervention change lasts for 30 days, we find that insurgents shift their attacks to the areas around Mosul while reducing the number of attacks in Baghdad. This displacement pattern is particularly pronounced for SAFs. For SAFs, insurgents appear to move their attacks to the Mosul area even with the intervention of 7 days, though the effect size is smaller. In short, the effects of airstrikes are not localized, but instead *ripple over long distances* as insurgents respond in different parts of the



(a) Increasing the expected number of airstrikes from 1 to 6 per day.



(b) Increasing the priority of Baghdad as focal point of airstrikes from $\alpha = 0$ to $\alpha = 3$.

Figure 7: Changes in Estimated Intensity of Insurgent Attacks when Increasing the Expected Number of Airstrikes (the first two rows) and when Shifting the Focal Point of Airstrikes to Baghdad (the bottom two rows). Insurgent attacks are measured using Improvised Explosive Devices (IEDs; the first and third rows) and small arms fire (SAFs; second and fourth rows) with the varying number of intervention duration, $M = 1, 3, 7, 30$ days (columns). The number shown below each map represents the estimated change in the total number of attacks per day over the entire country, whereas the legend represents the difference in estimated intensities.

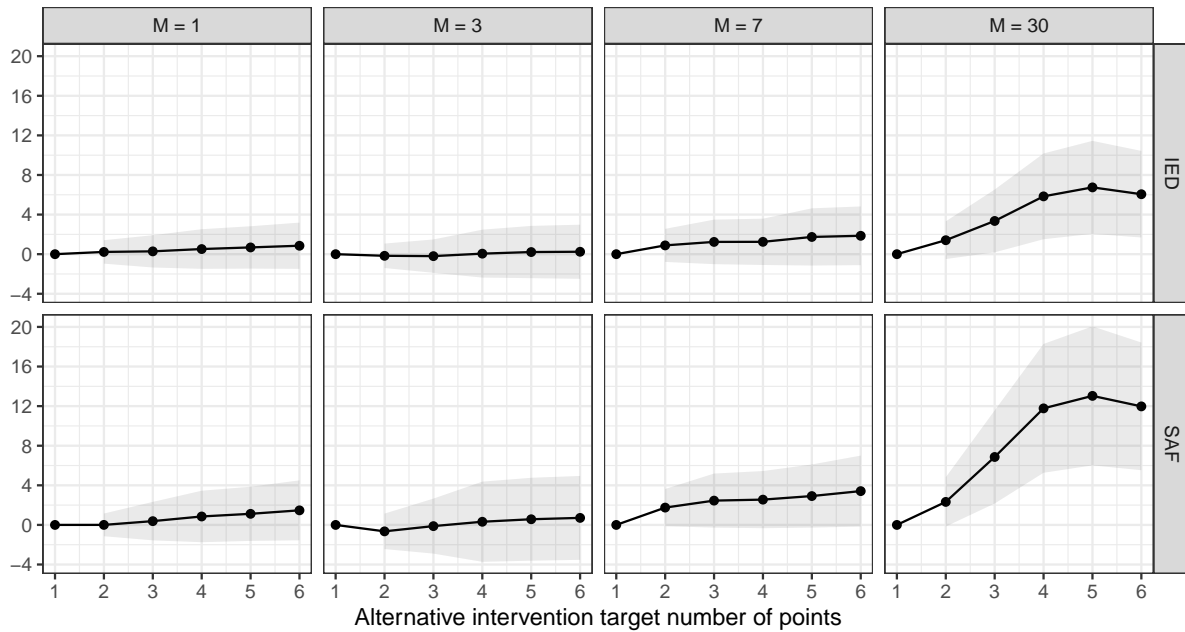
country. Unlike existing approaches which focus on the effect of an intervention in the nearby area (e.g. [Schutte and Donnay, 2014](#)), our approach captures this often-considerable displacement of violence and thus helps avoid mistaken inferences about the causal effects of airstrikes.

Figure 8a shows the changes in the estimated average number of insurgent attacks *in Baghdad* as the expected number of airstrikes increases from 1 to 2, 3, . . . , 6 airstrikes per day in the entire country (horizontal axis). We also vary the duration of intervention from $M = 1$ day to $M = 30$ days (columns). Both the point estimate (solid lines) and 95% CIs (grey bands) are shown. Consistent with Figure 7, we find that increasing the number of airstrikes leads to a greater number of attacks when the duration of intervention is 7 or 30 days. These effects appear are smaller when the intervention is much shorter. The patterns are similar for both IEDs and SAFs.

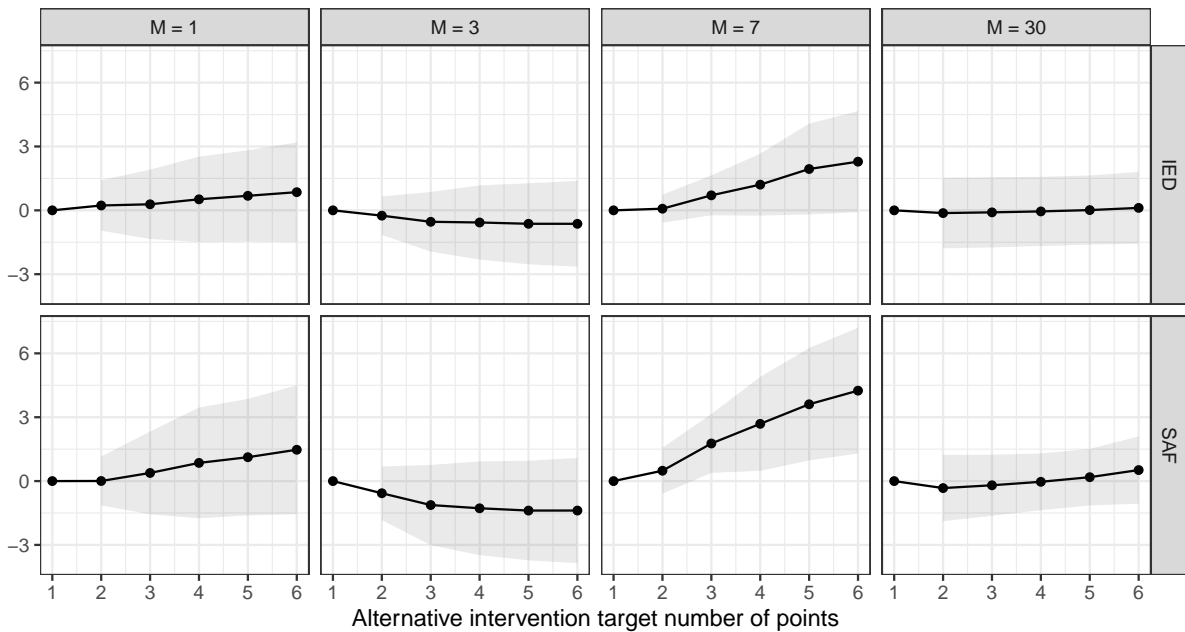
Figure 8b shows the change in the estimated number of IEDs and SAFs attacks in Baghdad when increasing the number of airstrikes M days before, while the expected number of airstrikes during the following $M - 1$ days equals one per day. We find that all estimated lagged effects for $M = 3$ are negative, whereas the estimated lagged effects for $M = 7$ are positive. This suggests that increasing the number of airstrikes may reduce insurgent violence in a short term while leading to an increase in a longer term. Appendix G.2 presents the effect estimates and 95% CIs for various interventions and outcomes.

We interpret these localized effects around Baghdad as consistent with prior claims (e.g. [Hashim, 2011](#)) that Sunni insurgents were sufficiently organized to shift their attacks to new fronts in response to American airstrikes. That is, while heavy bombardment in Baghdad might suppress insurgent attacks locally, we observe a net increase in overall violence as insurgent commanders displace their violence to new locations such as Mosul that are experiencing less airstrikes. This displacement effect underscores the danger of adopting too-narrow frameworks for casual estimation that miss spillover and other spatial knock-on effects.

The validity of our results hinges on the reliability of our causal assumptions, of which the uncon-



(a) Estimated effect of increasing the expected number of airstrikes for M days



(b) Estimated effect of increasing the expected number of airstrikes M days ago

Figure 8: Changes in the Estimated Number of Insurgency Attacks *in Baghdad* when Increasing the Expected Number of Airstrikes (a) for M days, and (b) M days ago. In each plot, the horizontal axis represents the expected number of airstrike per day under the alternative intervention. The vertical axis represents a change in the estimated average number of insurgency attacks in Baghdad for IEDs (first & third row) and SAFs (second & fourth row) when number of airstrikes per day increases from 1 to the value on the horizontal axis. Each column shows different (a) duration or (b) lag length of intervention, $M = 1, 3, 7, 30$ days. 95% confidence intervals are shown as grey bands.

foundedness assumption Assumption 1 is perhaps the strongest. We evaluate the robustness of our results to violations of this assumption using the sensitivity analysis framework developed in Section 5. We investigate the sensitivity of estimated effects for a change in intervention that corresponds to dosage or increased focus in Baghdad, for all values of M , for both SAF and IED outcomes, and for effects in the whole country and in Baghdad only. We find that the estimated effects are robust up to the ratio between the misspecified and the true propensity score (Γ) being bounded by 1.12. The small value of Γ indicates that our causal analysis may be sensitive to violations of the unconfoundedness assumption. As discussed before, however, this sensitivity is partially due to the inherently large uncertainty in estimating the point process intensity functions of the propensity scores from sparse data.

8 Concluding Remarks

In this paper, we provide a framework for causal inference with spatio-temporal point process treatments and outcomes. We demonstrate the flexibility of this proposed methodology by applying it to the estimation of airstrike effects on insurgent violence in Iraq. Our central idea is to use a stochastic intervention that represents a distribution of treatments rather than the standard causal inference approach that estimates the average potential outcomes under some fixed treatment values. A key advantage of our approach is its flexibility: it permits unstructured patterns of both spatial spillover and temporal carryover effects. This flexibility is crucial since for many spatio-temporal causal inference problems, including our own application, little is known about how the treatments in one area affect the outcomes in other areas across different time periods.

The estimands and methodology presented in this paper can be applied in a number of settings to estimate the effect of a particular stochastic intervention strategy. There are several considerations that may be useful when defining a stochastic intervention of interest. First, the choice of intervention should be guided by pressing policy questions or important academic debates where undetected spillover might frustrate traditional methods of causal inference. Second, stochastic interventions should satisfy the

overlap assumption (Assumption 2). Researchers should not define a stochastic intervention that generates treatment patterns that appear to be far different from those of the observed treatment events. In our application, we achieve this by constructing the stochastic interventions based on the estimated density of point patterns obtained from the past data and the observed number of airstrikes per day.

The proposed framework can also be applied to other high-dimensional, and possibly unstructured, treatments. The standard approach to causal inference, which estimates the causal effects of fixed treatment values, does not perform well in such settings. Indeed, the sparsity of observed treatment patterns alone makes it difficult to satisfy the required overlap assumption (Imai and Jiang, 2019). We believe that the stochastic intervention approach proposed here offers an effective solution to a broad class of causal inference problems.

Future research should further develop the methodology for stochastic interventions. In particular, it is important to consider an improved weighting method that explicitly targets covariate balance. This might be challenging in the spatiotemporal setting where the notion of covariate balance is not yet well understood. Finally, it is crucial to extend the stochastic intervention framework to adaptive strategies over multiple time periods that might be more reflective of realistic assignments.

References

- Abramson, I. S. (1982). On bandwidth variation in kernel estimates—a square root law. *The Annals of Statistics* 1217–1223.
- Aronow, P. M., Samii, C., and Wang, Y. (2019). Design-based inference for spatial experiments with interference. *Annual Summer Meeting of the Society for Political Methodology*.
- Baddeley, A., Rubak, E., and Turner, R. (2015). *Spatial Point Patterns: Methodology and Applications with R*. Chapman and Hall/CRC Press, London.
- Basse, G. and Airoidi, E. M. (2018). Limitations of design-based causal inference and a/b testing under arbitrary and network interference. *Sociological Methodology* **48**, 1, 136–151.

- Bojinov, I. and Shephard, N. (2019). Time series experiments and causal estimands: Exact randomization tests and trading. *Journal of the American Statistical Association* Forthcoming.
- Charnes, A. and Cooper, W. W. (1962). Programming with linear fractional functionals. *Naval Research logistics quarterly* **9**, 3-4, 181–186.
- Chow, Y. S. (1965). Local convergence of martingales and the law of large numbers. *The Annals of Mathematical Statistics* **36**, 2, 552–558.
- Cole, S. R., Edwards, J. K., Breskin, A., and Hudgens, M. G. (2021). Comparing parametric, non-parametric, and semiparametric estimators: The weibull trials. *American Journal of Epidemiology* .
- Crimaldi, I. and Pratelli, L. (2005). Convergence results for multivariate martingales. *Stochastic processes and their applications* **115**, 4, 571–577.
- Csörgö, M. (1968). On the Strong Law of Large Numbers and the Central Limit Theorem for Martingales. *Transactions of the American Mathematical Society* **131**, 1, 259–275.
- Dell, M. and Querubin, P. (2018). Nation building through foreign intervention: Evidence from discontinuities in military strategies. *Quarterly Journal of Economics* **133**, 2, 701–764.
- Díaz, I. and Hejazi, N. (2019). Causal mediation analysis for stochastic interventions. *arXiv preprint arXiv:1901.02776* .
- Díaz Muñoz, I. and van der Laan, M. (2012). Population Intervention Causal Effects Based on Stochastic Interventions. *Biometrics* **68**, 541–549.
- Fisher, R. A. (1935). *The Design of Experiments*. Oliver and Boyd, London.
- Gill, R. D. and Robins, J. M. (2001). Causal inference for longitudinal data: The continuous case. *Annals of Statistics* **29**, 6, 1785–1811.

- Hashim, A. (2011). *Insurgency and Counter-Insurgency in Iraq*. Cornell University Press, Ithaca.
- Hirano, K., Imbens, G. W., and Ridder, G. (2003). Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica* **71**, 4, 1161–1189.
- Hudgens, M. G. and Halloran, M. E. (2008). Toward Causal Inference With Interference. *Journal of the American Statistical Association* **103**, 482, 832–842.
- Imai, K. and Jiang, Z. (2019). Comment on “The Blessings of Multiple Causes” by Wang and Blei. *Journal of the American Statistical Association* **114**, 528, 1605–1610.
- Imai, K., Jiang, Z., and Malai, A. (2021). Causal inference with interference and noncompliance in two-stage randomized experiments. *Journal of the American Statistical Association* **116**, 534, 632–644.
- Kennedy, E. H. (2019). Nonparametric causal effects based on incremental propensity score interventions. *Journal of the American Statistical Association* **114**, 526, 645–656.
- Kocher, M., Pepinsky, T., and Kalyvas, S. (2011). Aerial bombing and counterinsurgency in the vietnam war. *American Journal of Political Science* **55**, 2, 201–218.
- Küchler, U., Sørensen, M., *et al.* (1999). A note on limit theorems for multivariate martingales. *Bernoulli* **5**, 3, 483–493.
- Liu, L., Hudgens, M. G., and Becker-Dreps, S. (2016). On inverse probability-weighted estimators in the presence of interference. *Biometrika* **103**, 4, 829–842.
- Lok, J. J. (2016). Defining and estimating causal direct and indirect effects when setting the mediator to specific values is not feasible. *Statistics in Medicine* **35**, 2, 4008–4020.
- Luo, X., Small, D. S., Li, C.-S. R., and Rosenbaum, P. R. (2012). Inference with interference between units in an fmri experiment of motor inhibition. *Journal of the American Statistical Association* **107**, 498, 530–541.

- Lyall, J. (2019a). Bombing to lose? airpower, civilian casualties, and the dynamics of violence in counterinsurgency wars. *Unpublished Paper* .
- Lyall, J. (2019b). Civilian casualties, humanitarian aid, and insurgent violence in civil wars. *International Organization* **73**, 4, 901–926.
- Mir, A. and Moore, D. (2019a). Drones, surveillance, and violence: Theory and evidence from a us drone program. *International Studies Quarterly* **63**, 4, 846–862.
- Mir, A. and Moore, D. (2019b). Drones, surveillance, and violence: Theory and evidence from a us drone program. *International Studies Quarterly* **63**, 4, 846–862.
- Neyman, J. (1923). On the application of probability theory to agricultural experiments: Essay on principles, section 9. (translated in 1990). *Statistical Science* **5**, 465–480.
- Papadogeorgou, G., Mealli, F., and Zigler, C. M. (2019). Causal inference with interfering units for cluster and population level treatment allocation programs. *Biometrics* **75**, 3, 778–787.
- Pearl, J. (2000). *Causality: Models, reasoning, and inference*. Cambridge: Cambridge University Press.
- Rigterink, A. (2021). The wane of command: Evidence on drone strikes and control within terrorist organizations. *American Political Science Review* **115**, 1, 31–50.
- Robins, J. (1986). A new approach to causal inference in mortality studies with a sustained exposure period—application to control of the healthy worker survivor effect. *Mathematical modelling* **7**, 9-12, 1393–1512.
- Robins, J. M. (1997). *Latent Variable Modeling and Applications to Causality*, vol. 120 of *Lecture Notes in Statistics*, chap. Causal Inference from Complex Longitudinal Data, 69–117. Springer Verlag, New York.
- Robins, J. M. (1999). Association, causation, and marginal structural models. *Synthese* 151–179.

- Robins, J. M., Hernán, M. A., and Brumback, B. (2000). Marginal structural models and causal inference in epidemiology. *Epidemiology* **11**, 5, 550–560.
- Rosenbaum, P. R. (2002). *Observational studies*. Springer.
- Rosenbaum, P. R. and Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika* **70**, 1, 41–55.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and non-randomized studies. *Journal of Educational Psychology* **66**, 688–701.
- Sävje, F., Aronow, P. M., and Hudgens, M. G. (2019). Average treatment effects in the presence of unknown interference. arXiv:1711.06399.
- Schutte, S. and Donnay, K. (2014). Matched wake analysis: Finding causal relationships in spatiotemporal event data. *Political Geography* **41**, 1–10.
- Scott, D. W. (1992). *Multivariate Density Estimation*. New York : Wiley.
- Serfling, R. J. (1980). *Approximation Theorems of Mathematical Statistics*. New York: Wiley.
- Sobel, M. E. and Lindquist, M. A. (2014). Causal inference for fmri time series data with systematic errors of measurement in a balanced on/off study of social evaluative threat. *Journal of the American Statistical Association* **109**, 507, 967–976.
- Stout, W. F. (1974). *Almost sure convergence*, vol. 24. Academic press.
- Tchetgen Tchetgen, E. J., Fulcher, I., and Shpitser, I. (2017). Auto-g-computation of causal effects on a network. *arXiv preprint arXiv:1709.01577* .
- Van der Vaart, A. W. (1998). *Asymptotic statistics*. Cambridge university press.
- van der Vaart, A. W. (2010). Time Series. *VU University Amsterdam, lecture notes* .

Young, J. G., Hernán, M. A., and Robins, J. M. (2014). Identification, estimation and approximation of risk under interventions that depend on the natural value of treatment using observational data. *Epidemiologic methods* **3**, 1, 1–19.

Zeng, S., Li, F., Hu, L., and Li, F. (2021). Propensity score weighting analysis for survival outcomes using pseudo observations. *arXiv preprint arXiv:2103.00605* .

Zhang, M., Joffe, M. M., and Small, D. S. (2011). Causal inference for continuous-time processes when covariates are observed only at discrete times. *Annals of Statistics* **39**, 1, 131–173.

SUPPLEMENTARY APPENDIX FOR “CAUSAL INFERENCE WITH SPATIO-TEMPORAL DATA”

Table of Contents

A	Notation	3
B	Theoretical Proofs	3
B.1	Regularity conditions	3
B.2	Proofs: The propensity score as a balancing score	6
B.3	Proofs: Asymptotic normality based on the true propensity score	7
B.4	Proofs: Asymptotic normality based on the estimated propensity score	14
B.5	Asymptotics for an increasing number of independent regions	34
C	The Hájek Estimator	38
D	Sensitivity analysis	39
D.1	For the IPW estimator	39
D.2	For the Hájek estimator	40
E	Additional Simulation Results on the Iraq-based scenario	41
E.1	Asymptotic Variance and Bound, and Estimated Variance Bound	41
E.2	Coverage of the Confidence Intervals for the IPW and Hájek Estimators	42
E.3	Uncertainty Estimates	45
E.4	Covariate Balance	47
F	Additional simulations on a square geometry	47

F.1	The Simulation Design	47
F.2	Simulation Results	49
G	Additional Empirical Results	56
G.1	Visualization	56
G.2	Empirical Results	59
G.3	Single time point adaptive interventions	59

A Notation

Table A.1: Notation.

Paths	\overline{W}_t	Treatments over the time periods $1, \dots, t$
	\overline{w}_t	Realized treatment assignments for time periods $1, \dots, t$
	$\overline{\mathcal{Y}}_t$	Collection of all potential outcomes for time periods $1, \dots, t$
	\overline{Y}_t	Observed outcomes for time periods $1, \dots, t$
Intervention	M	The number of time periods over which we intervene
	h	Poisson point process intensity defining the stochastic intervention
Estimands	N_t, N	Expected number of outcome-active locations during time period t for an intervention over M time periods, and their average over time
	τ_t^M, τ^M	Expected change in the number of outcome-active locations comparing two interventions for time period t and their average over time
Estimators	\widehat{Y}_t^M	Estimated continuous surface the integral of which is used for calculating \widehat{N}_t
	$\widehat{N}_t, \widehat{N}$	Estimated expected number of points during time period t for an intervention taking place over the preceding M time periods, and their average over time
	$\widehat{\tau}_t^M, \widehat{\tau}^M$	Estimated expected change in the number of outcome-active locations for time period t comparing two interventions, and their average over time
Arguments	B	The set over which the number of outcome-active locations are counted

B Theoretical Proofs

B.1 Regularity conditions

For $\epsilon > 0$, we use $\mathcal{N}_\epsilon(A)$ to denote the ϵ -neighborhood of a set A : $\mathcal{N}_\epsilon(A) = \{\omega \in \Omega : \text{there exists } a \in A \text{ with } \text{dist}(\omega, a) < \epsilon\}$. Also, we use ∂B to denote the boundary of B , formally defined as the set of points for which an open ball of any size centered at them includes points both in and outside B , i.e., $\partial B = \{s \in \Omega \text{ such that, for every } \epsilon > 0, \text{ there exist } s_1, s_2 \in \mathcal{N}_\epsilon(s) \text{ for which } s_1 \in B \text{ and } s_2 \notin B\}$.

Regularity conditions for asymptotic results when using the true or estimated propensity score

The following assumption includes regularity conditions which are used to show asymptotic normality of the estimator based on the true or estimated propensity score:

Assumption A.1. *The following three conditions hold.*

- (a) *There exists $\delta_Y > 0$ such that $|S_{Y_t(\bar{w}_t)}| < \delta_Y$ for all $t \in \mathcal{T}$ and $\bar{w}_t \in \mathcal{W}^T$.*
- (b) *Let $v_t = \text{Var} \left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)} N_B(Y_t) \mid \bar{H}_{t-M}^* \right]$ for $t \geq M$. Then, there exists $v \in \mathbb{R}^+$ such that $(T - M + 1)^{-1} \sum_{t=M}^T v_t \xrightarrow{p} v$ as $T \rightarrow \infty$.*
- (c) *There exists $\delta_B > 0$ and $Q^* \in (1/2, 1)$ such that*

$$P \left(\sum_{t=M}^T I \left(\exists s \in S_{Y_t} \cap \mathcal{N}_{\delta_B}(\partial B) \right) > T^{1-Q^*} \right) \rightarrow 0, \text{ as } T \rightarrow \infty.$$

Assumption A.1(a) states that there is an upper limit on the number of outcome-active locations at any time period and under any treatment path. In our application, it is reasonable to assume that the number of insurgent attacks occurring during any day is bounded. In Assumption A.1(b), \bar{H}_t^* represents the expanded history preceding W_{t+1} , including previous treatments, all potential outcomes, and all potential confounders. Given the assumptions of bounded relative positivity and bounded number of outcome-active locations, Assumption A.1(b) is a weak condition, as it states that the average of bounded quantities converges. Lastly, Assumption A.1(c) states that the probability that we observe more than T^{1-Q^*} time periods with outcome-active locations within a δ_B -neighborhood of B 's boundary goes to zero as the number of observed time periods increases. Since the size of the boundary's neighborhood can be arbitrarily small, this assumption is also reasonable. Informally, Assumption A.1(c) would be violated in our study if insurgent attacks occurred at the *boundary* of region B more often than during \sqrt{T} time periods. As long as the regions B are decided upon substantive interest, we would expect this assumption to be satisfied. Alternatively, regions B can be defined by avoiding setting the region's boundary at observed outcome-active locations.

Regularity conditions for asymptotic results when using the estimated propensity score

Next, we formalize the regularity conditions on the propensity score model. These conditions are used for establishing the asymptotic normality of the estimator based on the estimated propensity score.

Assumption A.2. *Assume that the parametric form of the propensity score indexed by γ , $f(W_t = w_t \mid \bar{H}_{t-1}; \gamma)$, is correctly specified and differentiable with respect to γ , and let $\psi(w_t, \bar{h}_{t-1}; \gamma) = \frac{\partial}{\partial \gamma} \log f(W_t = w_t \mid \bar{H}_{t-1} = \bar{h}_{t-1}; \gamma)$ be twice continuously differentiable score functions. Let γ_0 denote the true values of the parameters, where γ_0 is in an open subset of the Euclidean space. Denote*

$\mathcal{F}_t = \overline{H}_{t-M+1}^* = \{\overline{W}_{t-M+1}, \overline{Y}_T, \overline{X}_T\}$, as in the proof of Theorem 1. We assume that the following conditions hold:

1. (a) $E_{\gamma_0} \left[\|\psi(W_t, \overline{H}_{t-1}; \gamma_0)\|^2 \right] < \infty$,
(b) There exists a positive definite matrix V_{ps} such that

$$\frac{1}{T} \sum_{t=1}^T E_{\gamma_0} \left(\psi(W_t, \overline{H}_{t-1}; \gamma_0) \psi(W_t, \overline{H}_{t-1}; \gamma_0)^\top \mid \mathcal{F}_{t-1} \right) \xrightarrow{p} V_{ps}$$

$$(c) \frac{1}{T} \sum_{t=1}^T E_{\theta_0} \left[\|\psi(W_t, \overline{H}_{t-1}; \gamma_0)\|^2 I \left(\|\psi(W_t, \overline{H}_{t-1}; \gamma_0)\| > \epsilon \sqrt{T} \right) \mid \mathcal{F}_{t-1} \right] \xrightarrow{p} 0, \text{ for all } \epsilon > 0,$$

2. For all k, j , if we denote the k^{th} element of the $\psi(w_t, \overline{h}_{t-1}; \gamma)$ vector by $\psi_k(w_t, \overline{h}_{t-1}; \gamma)$ and $P_{kjt} = \frac{\partial}{\partial \gamma_j} \psi_k(W_t, \overline{H}_{t-1}; \gamma) \Big|_{\gamma_0}$, then $E_{\gamma_0} [|P_{kjt}|] < \infty$ and there exists $0 < r_{kj} \leq 2$ such that $\sum_{t=1}^T \frac{1}{t^{r_{kj}}} E_{\gamma_0} (|P_{kjt} - E_{\theta_0}(P_{kjt} \mid \mathcal{F}_{t-1})|^{r_{kj}} \mid \mathcal{F}_{t-1}) \xrightarrow{p} 0$
3. There exists an integrable function $\ddot{\psi}(w_t, \overline{h}_{t-1})$ such that $\ddot{\psi}(w_t, \overline{h}_{t-1})$ dominates the second partial derivatives of $\psi(w_t, \overline{h}_{t-1}; \gamma)$ in a neighborhood of γ_0 for all $(w_t, \overline{h}_{t-1})$.

Assumption A.3. Suppose that $\psi(w_t, \overline{h}_{t-1}; \gamma)$ are the score functions of a propensity score model that satisfies Assumption A.2 with true parameters γ_0 , and

$$s(\overline{h}_{t-1}, w_t, y_t; \gamma) = \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} \right] N_B(y_t) - N_{Bt}(F_h^M).$$

Then, the following conditions hold.

1. There exists $u \in \mathbb{R}^K$ such that

$$\frac{1}{T - M + 1} \sum_{t=M}^T E_{\theta_0} [s(\overline{H}_{t-1}, W_t, Y_t; \gamma_0) \psi(W_t, \overline{H}_{t-1}; \gamma_0) \mid \mathcal{F}_{t-1}] \xrightarrow{p} u,$$

2. If $P_{jt} = \frac{\partial}{\partial \gamma_j} s(\overline{H}_{t-1}, W_t, Y_t; \gamma) \Big|_{\gamma_0}$, where γ_j is the j^{th} entry of γ , then there exists $r_j \in (0, 2]$ such that

$$\sum_{t=1}^T \frac{1}{t^{r_j}} E_{\gamma_0} (|P_{jt} - E_{\theta_0}(P_{jt} \mid \mathcal{F}_{t-1})|^{r_j} \mid \mathcal{F}_{t-1}) \xrightarrow{p} 0.$$

Remark A.1. Given the previous assumptions, Assumption A.3 is quite weak. We look at the two parts separately:

1. For the k^{th} entry, we can write:

$$\begin{aligned}
& \frac{1}{T-M+1} \sum_{t=M}^T |E_{\theta_0}[s(\bar{H}_{t-1}, W_t, Y_t; \gamma_0) \psi_k(W_t, \bar{H}_{t-1}; \gamma_0) \mid \mathcal{F}_{t-1}]| \leq \\
& \leq \frac{1}{T-M+1} \sum_{t=M}^T \sqrt{E_{\theta_0}[s(\bar{H}_{t-1}, W_t, Y_t; \gamma_0)^2 \mid \mathcal{F}_{t-1}]} \sqrt{E_{\theta_0}[\psi_k(W_t, \bar{H}_{t-1}; \gamma_0)^2 \mid \mathcal{F}_{t-1}]} \\
& \hspace{20em} \text{(Cauchy-Schwarz)} \\
& \leq \frac{1}{2(T-M+1)} \sum_{t=M}^T \left(E_{\theta_0}[s(\bar{H}_{t-1}, W_t, Y_t; \gamma_0)^2 \mid \mathcal{F}_{t-1}] + E_{\theta_0}[\psi_k(W_t, \bar{H}_{t-1}; \gamma_0)^2 \mid \mathcal{F}_{t-1}] \right) \\
& \hspace{20em} (2ab \leq a^2 + b^2) \\
& \xrightarrow{p} \frac{1}{2} (v + [V_{ps}]_{kk}).
\end{aligned}$$

The proof that the first part converges to v will be shown in Equation (A.5), and the second part is based on Assumption A.2, where $[V_{ps}]_{kk}$ denotes the k^{th} diagonal entry of V_{ps} . Since the expression is already bounded at the limit, the assumption that it converges is reasonable. Furthermore, we have that $|u_k| \leq \frac{1}{2}(v + [V_{ps}]_{kk})$, where u_k is the k^{th} entry of u .

2. This assumption limits how much the derivative of $s(\bar{h}_{t-1}, w_t, y_t; \gamma)$ can vary around its conditional expectation. As we will see in Lemma A.2, this derivative can be re-written as a sum that involves three terms: the number of outcome active locations, the inverse probability ratios, and the score functions. The first two of these terms are bounded, and Assumption A.2 already controls how variable the score functions can be. Thus, this assumption is also reasonable.

B.2 Proofs: The propensity score as a balancing score

Proof of Proposition 1. Note that $f(W_t = w \mid e_t(w), \bar{H}_{t-1}) = f(W_t = w \mid \bar{H}_{t-1}) = e_t(w)$ since $e_t(w)$ is a function of \bar{H}_{t-1} . Therefore, it suffices to show that $f(W_t = w \mid e_t(w)) = e_t(w)$:

$$f(W_t = w \mid e_t(w)) = E[f(W_t = w \mid \bar{H}_{t-1}) \mid e_t(w)] = E[e_t(w) \mid e_t(w)] = e_t(w). \quad (\text{A.1})$$

□

Proof of Proposition 2.

$$f(W_t = w \mid \bar{\mathbf{W}}_{t-1}, \bar{\mathcal{Y}}_T, \bar{\mathcal{X}}_T)$$

$$\begin{aligned}
&= f(W_t = w \mid \bar{H}_{t-1}, \bar{\mathbf{W}}_{t-1}, \bar{\mathcal{Y}}_T, \bar{\mathcal{X}}_T) && \text{(Since } \bar{H}_{t-1} \subset \{\bar{\mathbf{W}}_{t-1}, \bar{\mathcal{Y}}_T, \bar{\mathcal{X}}_T\}) \\
&= f(W_t = w \mid \bar{H}_{t-1}) && \text{(From Assumption 1)} \\
&= e_t(w) \\
&= f(W_t = w \mid e_t(w)) && \text{(From (A.1))}
\end{aligned}$$

□

B.3 Proofs: Asymptotic normality based on the true propensity score

Proof of Theorem 1. Note that the collection of variables temporally precedent to treatment at time period t is the expanded history \bar{H}_{t-1}^* , defined in Assumption A.1. The expanded history \bar{H}_{t-1}^* is a filtration generated by the collection of potential confounders $\bar{\mathcal{X}}_T$, the collection of potential outcomes $\bar{\mathcal{Y}}_T$, and the previous treatments, and satisfies $\bar{H}_{t-1}^* \subset \bar{H}_t^*$.

Let $err_t = \hat{N}_{Bt}(F_h) - N_{Bt}(F_h)$ be the estimation error for time period t and lag M . We will decompose err_t in two components, one corresponding to the error due to the treatment assignment (A_{1t}), and the other corresponding to the error due to spatial smoothing (A_{2t}). Since the bandwidth parameter of the kernel depends on T , we write K_{b_T} instead of K_b . Specifically,

$$\begin{aligned}
err_t &= \left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)} \right] \int_B \sum_{s \in S_{Y_t}} K_{b_T}(\omega, s) d\omega - N_{Bt}(F_h^M) \\
&= \underbrace{\left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)} \right] N_B(Y_t) - N_{Bt}(F_h^M)}_{A_{1t}} + \\
&\quad \underbrace{\left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)} \right] \left[\int_B \sum_{s \in S_{Y_t}} K_{b_T}(\omega, s) d\omega - N_B(Y_t) \right]}_{A_{2t}}. \tag{A.2}
\end{aligned}$$

We show that

1. $\sqrt{T} \left(\frac{1}{T-M+1} \sum_{t=M}^T A_{1t} \right)$ is asymptotically normal, and
2. $\sqrt{T} \left(\frac{1}{T-M+1} \sum_{t=M}^T A_{2t} \right)$ converges to zero in probability.

Asymptotic normality of the first error.

We use the central limit theorem for martingale difference series (Theorem 4.16 of [van der Vaart \(2010\)](#)) to establish the asymptotic normality of $(T - M + 1)^{-1} \sum_{t=M}^T A_{1t}$.

Claim. A_{1t} is a martingale difference series with respect to the filtration $\mathcal{F}_t = \bar{H}_{t-M+1}^*$.

To prove this, we show that $E(|A_{1t}|) < \infty$ and $E(A_{1t} | \mathcal{F}_{t-1}) = E(A_{1t} | \overline{H}_{t-M}^*) = 0$. For the first part, Assumptions 2 and Assumption A.1(a) imply that A_{1t} is bounded and hence $E[|A_{1t}|] < \infty$:

$$|A_{1t}| \leq \left| \prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)} N_B(Y_t) \right| + |N_{Bt}(F_h^M)| \leq \delta_W^M \delta_Y + \delta_Y \quad (\text{A.3})$$

For the second part, it suffices to show

$$E \left\{ \left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)} \right] N_B(Y_t) \mid \overline{H}_{t-M}^* \right\} = N_{Bt}(F_h^M),$$

where the expectation is taken with respect to the assignment of treatments $\mathbf{W}_{(t-M+1):t}$.

$$\begin{aligned} & E \left\{ \left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)} \right] N_B(Y_t) \mid \overline{H}_{t-M}^* \right\} \\ &= \int \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j)} \right] N_B \left(Y_t(\overline{\mathbf{W}}_{t-M}, \underbrace{w_{t-M+1}, \dots, w_t}_{\mathbf{w}_{(t-M+1):t}}) \right) \times \\ & \quad f(w_{t-M+1} \mid \overline{H}_{t-M}^*) f(w_{t-M+2} \mid \overline{H}_{t-M}^*, W_{t-M+1}) \cdots \times \\ & \quad f(w_t \mid \overline{H}_{t-M}^*, \mathbf{W}_{(t-M+1):(t-1)}) d\mathbf{w}_{(t-M+1):t} \\ &= \int \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j)} \right] N_B \left(Y_t(\overline{\mathbf{W}}_{t-M}, w_{t-M+1}, \dots, w_t) \right) \times \\ & \quad f(w_{t-M+1} \mid \overline{H}_{t-M}^*) f(w_{t-M+2} \mid \overline{H}_{t-M+1}^*) \cdots f(w_t \mid \overline{H}_{t-1}^*) d\mathbf{w}_{(t-M+1):t} \\ & \quad \quad \quad \text{(because } \overline{H}_{t'+1}^* = \overline{H}_{t'}^* \cup \{W_{t'+1}\}) \\ &= \int N_B \left(Y_t(\overline{\mathbf{W}}_{t-M}, w_{t-M+1}, \dots, w_t) \right) \left[\prod_{j=t-M+1}^t f_h(w_j) \right] d\mathbf{w}_{(t-M+1):t} \quad \text{(By Assumption 1)} \\ &= N_{Bt}(F_h^M). \end{aligned} \quad (\text{A.4})$$

This proves that A_{1t} is a martingale difference series with respect to filtration \mathcal{F}_{t-1} .

Claim. $(T - M + 1)^{-1} \sum_{t=M}^T E\{A_{1t}^2 I(|A_{1t}| > \epsilon \sqrt{T - M + 1}) \mid \mathcal{F}_{t-1}\} \xrightarrow{P} 0$ for every $\epsilon > 0$.

Let $\epsilon > 0$. Note that A_{1t} is bounded by $\delta_Y(\delta_W^M + 1)$ (see Equation (A.3)). Choose T_0 as

$$\begin{aligned} T_0 &= \operatorname{argmin}_{t \in \mathbb{N}^+} \{ \epsilon \sqrt{t - M + 1} > \delta_Y(\delta_W^M + 1) \} \\ &= \operatorname{argmin}_{t \in \mathbb{N}^+} \left\{ t > M - 1 + \left[\frac{\delta_Y(\delta_W^M + 1)}{\epsilon} \right]^2 \right\} \end{aligned}$$

$$= \left[M - 1 + \left[\frac{\delta_Y(\delta_W^M + 1)}{\epsilon} \right]^2 \right].$$

Then, for $T > T_0$, we have that $\epsilon\sqrt{T - M + 1} > \epsilon\sqrt{T_0 - M + 1} > \delta_Y(\delta_W^M + 1)$ which leads to $I(|A_{1t}| > \epsilon\sqrt{T - M + 1}) = 0$ and $E(A_{1t}^2 I(|A_{1t}| > \epsilon\sqrt{T - M + 1}) | \mathcal{F}_{t-1}) = 0$. This proves the claim.

We now combine the above claims to establish the asymptotic normality of the first error. Since A_{1t} has mean zero, $E(A_{1t}^2 | \mathcal{F}_{t-1}) = \text{Var}(A_{1t} | \mathcal{F}_{t-1})$, and since $N_{Bt}(F_h)$ is fixed,

$$\text{Var}(A_{1t} | \mathcal{F}_{t-1}) = \text{Var} \left(\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)} N_B(Y_t) | \bar{H}_{t-M}^* \right)$$

which yields

$$\frac{1}{T - M + 1} \sum_{t=M}^T E(A_{1t}^2 | \mathcal{F}_{t-1}) \xrightarrow{p} v, \quad (\text{A.5})$$

from Assumption A.1(b). Combining these results, using that $\sqrt{T}/\sqrt{T - M + 1} \rightarrow 1$ and Theorem 4.16 of [van der Vaart \(2010\)](#), we have the desired result,

$$\sqrt{T} \left(\frac{1}{T - M + 1} \sum_{t=M}^T A_{1t} \right) \xrightarrow{d} N(0, v).$$

Convergence to zero of the second error.

The second error compares the integral of the kernel-smoothed outcome surface over the region of interest B with the actual number of points within the set B . We show that as T goes to infinity, and since the bandwidth of the kernel converges to 0, the error due to kernel smoothing also goes to zero. Specifically, we will show that

$$\sqrt{T} \left(\frac{1}{T - M + 1} \sum_{t=M}^T A_{2t} \right) \xrightarrow{p} 0.$$

Let $c_t = \prod_{j=t-M+1}^t f_h(W_j)/e_j(W_j)$, and write

$$\left| \frac{1}{T - M + 1} \sum_{t=M}^T A_{2t} \right| = \left| \frac{1}{T - M + 1} \sum_{t=M}^T c_t \left[\int_B \sum_{s \in S_{Y_t}} K_{b_T}(\omega; s) d\omega - N_B(Y_t) \right] \right|.$$

Then,

$$\int_B \sum_{s \in S_{Y_t}} K_{b_T}(\omega; s) d\omega - N_B(Y_t)$$

$$\begin{aligned}
&= \sum_{s \in S_{Y_t} \cap B} \int_B K_{b_T}(\omega; s) d\omega + \sum_{s \in S_{Y_t} \cap B^c} \int_B K_{b_T}(\omega; s) d\omega - N_B(Y_t) \\
&= \sum_{s \in S_{Y_t} \cap B} \left[1 - \int_{B^c} K_{b_T}(\omega; s) d\omega \right] + \sum_{s \in S_{Y_t} \cap B^c} \int_B K_{b_T}(\omega; s) d\omega - N_B(Y_t) \\
&= \sum_{s \in S_{Y_t} \cap B^c} \int_B K_{b_T}(\omega; s) d\omega - \sum_{s \in S_{Y_t} \cap B} \int_{B^c} K_{b_T}(\omega; s) d\omega.
\end{aligned}$$

This shows that the error from smoothing the outcome surface at time t comes from (1) the kernel weight from points outside of B that falls within B , and (2) the kernel weight from points inside B that falls outside B . Using this, we write:

$$\begin{aligned}
&\left| \frac{1}{T-M+1} \sum_{t=M}^T A_{2t} \right| = \\
&\left| \frac{1}{T-M+1} \sum_{t=M}^T c_t \left[\sum_{s \in S_{Y_t} \cap B^c} \int_B K_{b_T}(\omega; s) d\omega - \sum_{s \in S_{Y_t} \cap B} \int_{B^c} K_{b_T}(\omega; s) d\omega \right] \right|.
\end{aligned}$$

Take $\epsilon > 0$, and $Q \in (1/2, Q^*)$ where Q^* is the one in Assumption A.1(c). Then, we will show that $P(T^Q \{ \left| \frac{1}{T-M+1} \sum_{t=M}^T A_{2t} \right| > \epsilon \}) \rightarrow 0$ as $T \rightarrow \infty$, which implies that the second error converges to zero faster than \sqrt{T} (since $Q > 1/2$).

$$\begin{aligned}
&P\left(T^Q \left\{ \left| \frac{1}{T-M+1} \sum_{t=M}^T A_{2t} \right| > \epsilon \right\}\right) \\
&= P\left(\left| \frac{1}{T-M+1} \sum_{t=M}^T c_t \left[\sum_{s \in S_{Y_t} \cap B^c} \int_B K_{b_T}(\omega; s) d\omega - \sum_{s \in S_{Y_t} \cap B} \int_{B^c} K_{b_T}(\omega; s) d\omega \right] \right| > \frac{\epsilon}{T^Q}\right) \\
&\leq P\left(\frac{1}{T-M+1} \sum_{t=M}^T c_t \sum_{s \in S_{Y_t} \cap B^c} \int_B K_{b_T}(\omega; s) d\omega > \frac{\epsilon}{2T^Q}\right) + \\
&\quad P\left(\frac{1}{T-M+1} \sum_{t=M}^T c_t \sum_{s \in S_{Y_t} \cap B} \int_{B^c} K_{b_T}(\omega; s) d\omega > \frac{\epsilon}{2T^Q}\right),
\end{aligned}$$

where the last equation holds because $|A - B| > \epsilon$ implies that at least one of $|A|, |B| > \epsilon/2$. Also, since all quantities are positive, we can drop the absolute value. Then, since $c_t \leq \delta_W^M$ from Assumption 2,

$$\begin{aligned}
&P\left(T^Q \left\{ \left| \frac{1}{T-M+1} \sum_{t=M}^T A_{2t} \right| > \epsilon \right\}\right) \\
&\leq P\left(\frac{1}{T-M+1} \sum_{t=M}^T \sum_{s \in S_{Y_t} \cap B^c} \int_B K_{b_T}(\omega; s) d\omega > \frac{\epsilon}{2T^Q \delta_W^M}\right) +
\end{aligned}$$

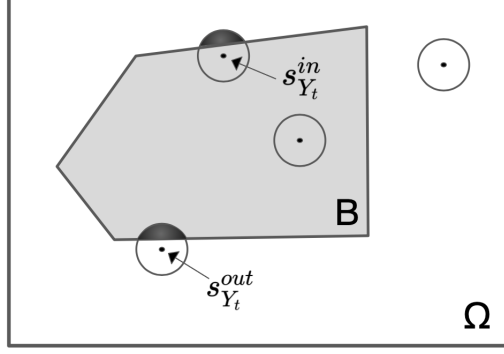


Figure A.1: Kernel-smoothed outcome surface, and points $s_{Y_t}^{in}, s_{Y_t}^{out}$ as the points closest to the boundary of B that lie within and outside B respectively. The amount of kernel weight falling within B from points outside of B is necessarily less or equal to the kernel weight from $s_{Y_t}^{out}$ (shaded), and similarly for $s_{Y_t}^{in}$.

$$P\left(\frac{1}{T-M+1} \sum_{t=M}^T \sum_{s \in S_{Y_t} \cap B} \int_{B^c} K_{b_T}(\omega; s) d\omega > \frac{\epsilon}{2T^Q \delta_W^M}\right).$$

Use $s_{Y_t}^{out}$ to denote the point in S_{Y_t} that lies outside B and is the closest to B : $s_{Y_t}^{out} = \{s \in S_{Y_t} \cap B^c : \text{dist}(s, B) = \min_{s' \in S_{Y_t} \cap B^c} \text{dist}(s', B)\}$. Similarly, $s_{Y_t}^{in}$ is the point in $S_{Y_t} \cap B$ that is closest to B^c . These points are shown graphically in Figure A.1. Because there are at most δ_Y outcome-active locations, from the definition of $s_{Y_t}^{in}, s_{Y_t}^{out}$, and because kernels are defined to be decreasing in distance, we have that

$$\begin{aligned} & P\left(T^Q \left\{ \left| \frac{1}{T-M+1} \sum_{t=M}^T A_{2t} \right| \right\} > \epsilon\right) \\ & \leq P\left(\frac{1}{T-M+1} \sum_{t=M}^T \int_B K_{b_T}(\omega; s_{Y_t}^{out}) d\omega > \frac{\epsilon}{2T^Q \delta_W^M \delta_Y}\right) \\ & \quad + P\left(\frac{1}{T-M+1} \sum_{t=M}^T \int_{B^c} K_{b_T}(\omega; s_{Y_t}^{in}) d\omega > \frac{\epsilon}{2T^Q \delta_W^M \delta_Y}\right) \\ & = \underbrace{P\left(\sum_{t=M}^T \int_B K_{b_T}(\omega; s_{Y_t}^{out}) d\omega > \frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y}\right)}_{B_1} \\ & \quad + \underbrace{P\left(\sum_{t=M}^T \int_{B^c} K_{b_T}(\omega; s_{Y_t}^{in}) d\omega > \frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y}\right)}_{B_2}. \end{aligned}$$

We show that B_1, B_2 converge to zero separately. Take B_1 :

$$\begin{aligned}
B_1 &= P\left(\sum_{t=M}^T \int_B K_{b_T}(\omega; s_{Y_t}^{out}) d\omega > \frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y} \left| \sum_{t=M}^T I(s_{Y_t}^{out} \in \mathcal{N}_{\delta_B}(\partial B)) > T^{1-Q^*} \right.\right) \\
&\quad \times P\left(\sum_{t=M}^T I(s_{Y_t}^{out} \in \mathcal{N}_{\delta_B}(\partial B)) > T^{1-Q^*}\right) \\
&+ P\left(\sum_{t=M}^T \int_B K_{b_T}(\omega; s_{Y_t}^{out}) d\omega > \frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y} \left| \sum_{t=M}^T I(s_{Y_t}^{out} \in \mathcal{N}_{\delta_B}(\partial B)) \leq T^{1-Q^*} \right.\right) \\
&\quad \times P\left(\sum_{t=M}^T I(s_{Y_t}^{out} \in \mathcal{N}_{\delta_B}(\partial B)) \leq T^{1-Q^*}\right)
\end{aligned}$$

From Assumption (c) we have that

$$\begin{aligned}
&P\left(\sum_{t=M}^T I(s_{Y_t}^{out} \in \mathcal{N}_{\delta_B}(\partial B)) > T^{1-Q^*}\right) \\
&\leq P\left(\sum_{t=M}^T I(\exists s \in S_{Y_t} \cap \mathcal{N}_{\delta_B}(\partial B)) > T^{1-Q^*}\right) \rightarrow 0,
\end{aligned}$$

and $\lim_{T \rightarrow \infty} B_1$ is equal to

$$\lim_{T \rightarrow \infty} P\left(\sum_{t=M}^T \int_B K_{b_T}(\omega; s_{Y_t}^{out}) d\omega > \frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y} \left| \sum_{t=M}^T I(s_{Y_t}^{out} \in \mathcal{N}_{\delta_B}(\partial B)) \leq T^{1-Q^*} \right.\right).$$

Studying the latter quantity, we have that

$$\begin{aligned}
&P\left(\sum_{t=M}^T \int_B K_{b_T}(\omega; s_{Y_t}^{out}) d\omega > \frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y} \left| \sum_{t=M}^T I(s_{Y_t}^{out} \in \mathcal{N}_{\delta_B}(\partial B)) \leq T^{1-Q^*} \right.\right) \\
&\leq P\left(\sum_{\substack{t=M \\ s_{Y_t}^{out} \notin \mathcal{N}_{\delta_B}(\partial B)}}^T \int_B K_{b_T}(\omega; s_{Y_t}^{out}) d\omega > \frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y} - T^{1-Q^*}\right) \\
&\leq P\left(\sum_{\substack{t=M \\ s_{Y_t}^{out} \notin \mathcal{N}_{\delta_B}(\partial B)}}^T \int_{\omega: \|\omega\| > \delta_B} K_{b_T}(\omega; \mathbf{0}) d\omega > \frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y} - T^{1-Q^*}\right) \\
&\leq P\left((T-M+1) \int_{\omega: \|\omega\| > \delta_B} K_{b_T}(\omega; \mathbf{0}) d\omega > \frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y} - T^{1-Q^*}\right) \\
&= I\left((T-M+1) \int_{\omega: \|\omega\| > \delta_B} K_{b_T}(\omega; \mathbf{0}) d\omega > \frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y} - T^{1-Q^*}\right) \tag{A.6}
\end{aligned}$$

where the first inequality follows from the fact that at most T^{1-Q^*} time periods had $s_{Y_t}^{out}$ within δ_B of set's B boundary, and $\int_B K_{b_T}(\omega; s_{Y_t}^{out}) \leq 1$ for those time periods. The second inequality follows from the fact that during the remaining time periods $s_{Y_t}^{out}$ was further than δ_B from B and $\int_B K_{b_T}(\omega; s_{Y_t}^{out}) \leq \int_{\omega: \|\omega - s_{Y_t}^{out}\| > \delta_B} K_{b_T}(\omega; s_{Y_t}^{out}) = \int_{\omega: \|\omega\| > \delta_B} K_{b_T}(\omega; \mathbf{0})$. The third inequality follows from not excluding the time periods with $s_{Y_t}^{out} \in \mathcal{N}_{\delta_B}(\partial B)$. Finally, the last equality holds because there is no uncertainty in the statement so the probability turns to an indicator.

Since $b_T \rightarrow 0$ as $T \rightarrow 0$, there exists $T_1 \in \mathbb{N}$ such that $b_T < \delta_B$ and $\int_{\omega: \|\omega\| > \delta_B} K_{b_T}(\omega; \mathbf{0}) d\omega = 0$ for all $T \geq T_1$. Also, since $\frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y} - T^{1-Q^*} \rightarrow \infty$, there exists $T_2 \in \mathbb{N}$ such that $\frac{\epsilon(T-M+1)}{2T^Q \delta_W^M \delta_Y} - T^{1-Q^*} > 1$ for all $T \geq T_2$. Then, for all $T \geq T_0 = \max\{T_1, T_2\}$ we have that the quantity in Equation (A.6) is equal to 0, showing that $\lim_{T \rightarrow \infty} B_1 = 0$. Similarly, we can show that $\lim_{T \rightarrow \infty} B_2 = 0$.

Combining all of these results we have that

$$P\left(T^Q \left\{ \left| \frac{1}{T-M+1} \sum_{t=M}^T A_{2t} \right| \right\} > \epsilon\right) \rightarrow 0,$$

as $T \rightarrow \infty$, establishing that the second error converges to zero faster than $1/\sqrt{T}$. \square

Proof of Lemma 1. Define $\Psi_t = [\widehat{N}_{Bt}(F_h^M)]^2 - v_t^*$. Then, Ψ_t is a martingale difference series with respect to $\mathcal{F}_t = \overline{H}_{t-M+1}$ since the following two hold: (1) $E(|\Psi_t|) < \infty$ since Ψ_t is bounded, and (2) $E(\Psi_t | \mathcal{F}_{t-1}) = E\left\{ [\widehat{N}_{Bt}(F_h^M)]^2 | \overline{H}_{t-M}^* \right\} - v_t^* = 0$. Also, since $\widehat{N}_{Bt}(F_h^M)$ is bounded we have that $\sum_{t=M}^{\infty} t^{-2} E(\Psi_t^2) < \infty$. From Theorem 1 in Csörgö (1968) we have that

$$\frac{1}{T-M+1} \sum_{t=M}^T \Psi_t = \frac{1}{T-M+1} \sum_{t=M}^T [\widehat{N}_{Bt}(F_h^M)]^2 - \frac{1}{T-M+1} \sum_{t=M}^T v_t^* \xrightarrow{p} 0.$$

\square

We use the results above to acquire asymptotic normality of the estimator for the causal effect, $\widehat{\tau}_B(F_{h_1}^M, F_{h_2}^M)$:

Theorem A.1. *Suppose that Assumptions 1 and 2 as well as the regularity conditions (Assumption A.1) hold. If the bandwidth $b_T \rightarrow 0$, then we have that*

$$\sqrt{T}(\widehat{\tau}_B(F_{h_1}^M, F_{h_2}^M) - \tau_B(F_{h_1}^M, F_{h_2}^M)) \xrightarrow{d} N(0, \eta),$$

as $T \rightarrow \infty$, for some $\eta > 0$. Finally, an upper bound of the asymptotic variance η can be consistently estimated by

$$\frac{1}{T-M+1} \sum_{t=M}^T [\widehat{\tau}_{Bt}(F_{h_1}^M, F_{h_2}^M)]^2 \xrightarrow{p} \eta^* \geq \eta.$$

Proof. In order to prove the asymptotic normality of $\widehat{\tau}_{Bt}(F_{h_1}^M, F_{h_2}^M)$ we will rely on results in the proof of Theorem 1 above. Take

$$\begin{aligned} \widehat{\tau}_{Bt}(F_{h_1}^M, F_{h_2}^M) - \tau_{Bt}(F_{h_1}^M, F_{h_2}^M) &= \\ & \left\{ \prod_{j=t-M+1}^t \frac{f_{h_2}(W_j)}{e_j(W_j)} - \prod_{j=t-M+1}^t \frac{f_{h_1}(W_j)}{e_j(W_j)} \right\} \int_B \sum_{s \in S_{Y_t}} K_{b_T}(\omega, s) d\omega - \tau_{Bt}(F_{h_1}^M, F_{h_2}^M) = \\ & \underbrace{\left\{ \prod_{j=t-M+1}^t \frac{f_{h_2}(W_j)}{e_j(W_j)} - \prod_{j=t-M+1}^t \frac{f_{h_1}(W_j)}{e_j(W_j)} \right\} N_B(Y_t) - \tau_{Bt}(F_{h_1}^M, F_{h_2}^M)}_{C_{1t}} + \\ & \underbrace{\left[\prod_{j=t-M+1}^t \frac{f_{h_2}(W_j)}{e_j(W_j)} \right] \left[\int_B \sum_{s \in S_{Y_t}} K_{b_T}(\omega, s) d\omega - N_B(Y_t) \right]}_{C_{2t}^2} - \\ & \underbrace{\left[\prod_{j=t-M+1}^t \frac{f_{h_1}(W_j)}{e_j(W_j)} \right] \left[\int_B \sum_{s \in S_{Y_t}} K_{b_T}(\omega, s) d\omega - N_B(Y_t) \right]}_{C_{2t}^1} \end{aligned}$$

Following steps identical to showing $\sqrt{T} \left[(T-M+1)^{-1} \sum_{t=M}^T A_{2t} \right] \xrightarrow{p} 0$ in the proof of Theorem 1, we can equivalently show that $\sqrt{T} \left[(T-M+1)^{-1} \sum_{t=M}^T C_{2t}^1 \right] \xrightarrow{p} 0$ and $\sqrt{T} \left[(T-M+1)^{-1} \sum_{t=M}^T C_{2t}^2 \right] \xrightarrow{p} 0$.

Therefore, all we need to show is that $\sqrt{T} \left[(T-M+1)^{-1} \sum_{t=M}^T C_{1t} \right] \xrightarrow{d} N(0, \eta)$. We will do so by showing again that C_{1t} is a martingale difference series with respect to the filtration \mathcal{F}_{t-1} :

1. Since $E(|A_{1t}|) < \infty$, from the triangular inequality we straightforwardly have that $E(|C_{1t}|) < \infty$.
2. Since $E(A_{1t} | \mathcal{F}_{t-1}) = 0$, we also have that $E(C_{1t} | \mathcal{F}_{t-1}) = 0$, from linearity of expectation.

Then, using the triangular inequality and Equation (A.3), we have that C_{1t} is bounded by $2\delta_Y(\delta_W^M + 1)$. Then, for $\epsilon > 0$, choosing $T_0 = \underset{t \in \mathbb{N}^+}{\operatorname{argmin}} \{ \epsilon \sqrt{t-M+1} > 2\delta_Y(\delta_W^M + 1) \}$ satisfies that, for $T > T_0$, $E(C_{1t}^2 I(|C_{1t}| > \epsilon \sqrt{T-M+1}) | \mathcal{F}_{t-1}) = 0$. Combining these results, we have that $\sqrt{T} [\widehat{\tau}_{Bt}(F_{h_1}^M, F_{h_2}^M) - \tau_{Bt}(F_{h_1}^M, F_{h_2}^M)] \rightarrow N(0, \eta)$.

To show $(T-M+1)^{-1} \sum_{t=M}^T \left\{ [\widehat{\tau}_{Bt}(F_{h_1}^M, F_{h_2}^M)]^2 - E\{[\tau_{Bt}(F_{h_1}^M, F_{h_2}^M)]^2 | \overline{H}_{t-M}^*\} \right\} \xrightarrow{p} 0$, the proof follows exactly the same way as the proof of Lemma 1 and is omitted here. \square

B.4 Proofs: Asymptotic normality based on the estimated propensity score

We will prove the asymptotic normality of the proposed estimators when the propensity score is estimated using a correctly specified parametric model. We extend Theorem 4.16 of [van der Vaart \(2010\)](#) to

multivariate martingale difference series. To our knowledge, this result is new even though the related results exist in the continuous time setting (Küchler *et al.*, 1999; Crimaldi and Pratelli, 2005). Under some additional assumptions on the martingale series, we show that the solution to the empirical estimating equation is also asymptotically normal. This result will be crucial in establishing the asymptotic normality of the maximum likelihood estimator for the propensity score model parameters. Finally, we combine these results and apply them to our specific context.

Theorem A.2 (Central limit theorem for multivariate martingale difference series). *Let $X_t = (X_{1t}, X_{2t}, \dots, X_{Kt})^\top$ be a multivariate martingale difference series with respect to the filtration \mathcal{F}_t in that $E[X_t | \mathcal{F}_{t-1}] = 0$ and $E[\|X_t\|] < \infty$, where $\|X_t\| = \sqrt{X_t^\top X_t} = \sqrt{\sum_{k=1}^K X_{kt}^2}$. Suppose that the following conditions hold.*

1. *There exists positive definite matrix $V \in \mathbb{R}^{K \times K}$ such that $\frac{1}{T} \sum_{t=1}^T E(X_t X_t^\top | \mathcal{F}_{t-1}) \xrightarrow{p} V$,*
2. $\frac{1}{T} \sum_{t=1}^T E\left[\|X_t\|^2 I\left(\|X_t\| > \epsilon\sqrt{T}\right) | \mathcal{F}_{t-1}\right] \xrightarrow{p} 0$, *for all $\epsilon > 0$.*

Then, we have,

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T X_t \xrightarrow{d} N(0, V).$$

Proof. We will use the Cramer-Wold device. We show that for every $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_K) \in \mathbb{R}^K$, it holds that $\alpha^\top \frac{1}{\sqrt{T}} \sum_{t=1}^T X_t \xrightarrow{d} \alpha^\top N(0, V)$. If this is true, it is implied that

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T X_t \xrightarrow{d} N(0, V).$$

Clearly, if α is the zero-vector, the result is trivial. So we focus on vectors α such that $\|\alpha\| \neq 0$.

Define $Y_t = \alpha^\top X_t$. First we show that Y_t is a martingale difference series with respect to \mathcal{F}_t :

$$E(|Y_t|) \leq \sum_{k=1}^K |\alpha_k| E(|X_{kt}|) \leq \sum_{k=1}^K |\alpha_k| E(\|X_t\|) < \infty, \quad \text{and}$$

$$E(Y_t | \mathcal{F}_{t-1}) = \sum_{k=1}^K \alpha_k E(X_{kt} | \mathcal{F}_{t-1}) = 0,$$

since X_t is a martingale difference series with respect to \mathcal{F}_t . So Y_t is also a martingale difference series with respect to \mathcal{F}_t . Next we will show that the conditions of Theorem 4.16 of van der Vaart (2010) hold

for Y_t . For Condition 2 we will use the fact that $\alpha\alpha^\top$ is a rank 1 symmetric matrix of dimension K with only non-zero eigenvalue equal to $\|\alpha\|^2$, and for that reason $X_t^\top\alpha\alpha^\top X_t = \|\alpha\|^2 X_t^\top X_t = \|\alpha\|^2 \|X_t\|^2$.

$$\begin{aligned} \text{Condition 1} \quad & \frac{1}{T} \sum_{t=1}^T E(Y_t^2 \mid \mathcal{F}_{t-1}) = \frac{1}{T} \sum_{t=1}^T E\left(\alpha^\top X_t X_t^\top \alpha \mid \mathcal{F}_{t-1}\right) \\ & = \alpha^\top \frac{1}{T} \sum_{t=1}^T E\left(X_t X_t^\top \mid \mathcal{F}_{t-1}\right) \alpha \xrightarrow{p} \alpha^\top V \alpha \\ & \hspace{15em} (\text{from the first assumption of Theorem A.2}) \end{aligned}$$

$$\begin{aligned} \text{Condition 2} \quad & E\left[Y_t^2 I(|Y_t| > \epsilon\sqrt{n}) \mid \mathcal{F}_{t-1}\right] = \\ & = E\left[X_t^\top \alpha \alpha^\top X_t I(X_t^\top \alpha \alpha^\top X_t > \epsilon^2 n) \mid \mathcal{F}_{t-1}\right] \\ & = E\left[\|\alpha\|^2 \|X_t\|^2 I(\|\alpha\|^2 \|X_t\|^2 > \epsilon^2 n) \mid \mathcal{F}_{t-1}\right] \\ & = \|\alpha\|^2 E\left[\|X_t\|^2 I\left(\|X_t\| > \frac{\epsilon}{\|\alpha\|} \sqrt{n}\right) \mid \mathcal{F}_{t-1}\right] \xrightarrow{p} 0 \\ & \hspace{15em} (\text{From the second condition of the Theorem for } \epsilon' = \epsilon/\|\alpha\|) \end{aligned}$$

Using Theorem 4.16 from [van der Vaart \(2010\)](#):

$$\alpha^\top \frac{1}{\sqrt{T}} \sum_{t=1}^T X_t = \sqrt{T} \frac{1}{T} \sum_{t=1}^T Y_t \xrightarrow{d} N(0, \alpha^\top V \alpha) \stackrel{d}{=} \alpha^\top N(0, V).$$

□

Now that a multivariate central limit theorem (CLT) for martingale difference series is established, we prove the next result which will be crucial in obtaining the asymptotic normality of estimators for the propensity score parameters. To our knowledge, this result is also new in martingale theory, but a related result in the iid setting is given as Theorem 5.21 of [Van der Vaart \(1998\)](#).

Theorem A.3 (Asymptotic normality of the solution to the estimating equation). *Let $\theta \rightarrow s(x, \theta) = (s_1(x, \theta), s_2(x, \theta), \dots, s_K(x, \theta))^\top \in \mathbb{R}^K$ be twice continuously differentiable with respect to $\theta = (\theta_1, \theta_2, \dots, \theta_K)^\top \in \Theta$, open subset of \mathbb{R}^K . Suppose that the following conditions hold.*

1. $s(X_t, \theta_0)$ satisfies the conditions of Theorem A.2 under θ_0 , in that there exists filtration \mathcal{F}_t such that

(a) $E_{\theta_0}[s(X_t, \theta_0) \mid \mathcal{F}_{t-1}] = 0$ and $E_{\theta_0}[\|s(X_t, \theta_0)\|] < \infty$ (and therefore it is a martingale difference series),

(b) $\exists V \in \mathbb{R}^{K \times K}$ positive definite such that $\frac{1}{T} \sum_{t=1}^T E_{\theta_0}\left(s(X_t, \theta_0) s(X_t, \theta_0)^\top \mid \mathcal{F}_{t-1}\right) \xrightarrow{p} V$, and

- (c) $\frac{1}{T} \sum_{t=1}^T E_{\theta_0} \left[\|s(X_t, \theta_0)\|^2 I \left(\|s(X_t, \theta_0)\| > \epsilon \sqrt{T} \right) \mid \mathcal{F}_{t-1} \right] \xrightarrow{p} 0$, for all $\epsilon > 0$,
2. $\frac{1}{T} \sum_{t=1}^T E_{\theta_0} \left(\left. \frac{\partial}{\partial \theta^T} s(X_t, \theta) \right|_{\theta_0} \mid \mathcal{F}_{t-1} \right) \xrightarrow{p} V_d$, for $V_d \in \mathbb{R}^{K \times K}$ invertible,
3. for all k, j , if we denote $P_{kjt} = \left. \frac{\partial}{\partial \theta_j} s_k(X_t, \theta) \right|_{\theta_0}$, we have that $E_{\theta_0}[|P_{kjt}|] < \infty$, and there exists $0 < r_{kj} \leq 2$ such that $\sum_{t=1}^T \frac{1}{t^{r_{kj}}} E_{\theta_0} (|P_{kjt} - E_{\theta_0}[P_{kjt} \mid \mathcal{F}_{t-1}]|^{r_{kj}} \mid \mathcal{F}_{t-1}) \xrightarrow{p} 0$,
4. there exists an integrable function $\ddot{\psi}(x)$ such that $\ddot{\psi}(x)$ dominates the second partial derivatives of $s_k(x, \theta)$ in a neighborhood of θ_0 for all x , and $k = 1, 2, \dots, K$.

If $\Psi_T(\theta) = \frac{1}{T} \sum_{t=1}^T s(X_t, \theta)$, and the solution to $\Psi_T(\theta) = 0$, $\hat{\theta}_T$, is consistent for θ_0 , then

$$\sqrt{T} (\hat{\theta}_T - \theta_0) \xrightarrow{d} N(0, V_d^{-1} V (V_d^{-1})^\top).$$

Proof. We extend the proof of Theorem 5.41 of [Van der Vaart \(1998\)](#) from the iid to the time series setting. Since the conditions of Theorem [A.2](#) are satisfied under θ_0 , we have that

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T s(X_t, \theta_0) = \sqrt{T} \Psi_T(\theta_0) \rightarrow N(0, V).$$

We will use the Taylor expansion for the vector valued $\Psi_T(\hat{\theta})$ around $\theta_0 = (\theta_{01}, \theta_{02}, \dots, \theta_{0K})^T$. To do so, we define the matrix $\dot{\Psi}_T(\theta) \in \mathbb{R}^{K \times K}$ and array $\ddot{\Psi}_T(\theta) \in \mathbb{R}^{K \times K \times K}$ of first and second derivatives as

$$\begin{aligned} \left[\dot{\Psi}_T(\theta) \right]_{kj} &= \left. \frac{\partial}{\partial \theta_j} \Psi_{kT}(\theta) \right|_{\theta} = \frac{1}{T} \sum_{t=1}^T \left. \frac{\partial}{\partial \theta_j} s_k(X_t, \theta) \right|_{\theta} && \text{and} \\ \left[\ddot{\Psi}_T(\theta) \right]_{kji} &= \left. \frac{\partial^2}{\partial \theta_j \partial \theta_i} \Psi_{kT}(\theta) \right|_{\theta} = \frac{1}{T} \sum_{t=1}^T \left. \frac{\partial^2}{\partial \theta_j \partial \theta_i} s_k(X_t, \theta) \right|_{\theta}, \end{aligned}$$

for $i, j, k = 1, 2, \dots, q$, where Ψ_{kT} is the k^{th} element of the Ψ_T vector. Then, we can write the Taylor expansion as

$$\Psi_T(\hat{\theta}) = \Psi_T(\theta_0) + \dot{\Psi}_T(\theta_0)(\hat{\theta}_T - \theta_0) + \ddot{\Psi}_T(\theta^*)(\hat{\theta}_T - \theta_0, \hat{\theta}_T - \theta_0), \quad (\text{A.7})$$

where θ^* is between $\widehat{\theta}_T$ and θ_0 , and $\ddot{\Psi}_T(\theta^*)(\widehat{\theta}_T - \theta_0, \widehat{\theta}_T - \theta_0)$ is a vector of length K with k^{th} entry

$$\sum_{j,i=1}^K \left[\ddot{\Psi}_T(\theta^*) \right]_{kji} (\widehat{\theta}_{Tj} - \theta_{0j})(\widehat{\theta}_{Ti} - \theta_{0i}),$$

and $\widehat{\theta}_{Ti}$ is the i^{th} entry of $\widehat{\theta}$. Therefore, we can write $\ddot{\Psi}_T(\theta^*)(\widehat{\theta}_T - \theta_0, \widehat{\theta}_T - \theta_0)$ as $A_T(\widehat{\theta}_T - \theta_0)$ where A_T is the $K \times K$ matrix which is the result of multiplying the tensor $\ddot{\Psi}_T(\theta^*)$ with the vector $\widehat{\theta}_T - \theta_0$ along the second mode, and it has (k, i) entry equal to

$$[A_T]_{ki} = \sum_{j=1}^K \left[\ddot{\Psi}_T(\theta^*) \right]_{kji} (\widehat{\theta}_{Tj} - \theta_{0j})$$

For notational simplicity, we do not include θ^* and $\widehat{\theta}_T - \theta_0$ in the notation of A_T . Since $\widehat{\theta}$ is the solution to $\Psi_T(\theta) = 0$, and based on the above, we can re-write Equation (A.7) as

$$\begin{aligned} 0 &= \Psi_T(\theta_0) + \dot{\Psi}_T(\theta_0)(\widehat{\theta}_T - \theta_0) + A_T(\widehat{\theta}_T - \theta_0) \\ \implies -\sqrt{T}\Psi_T(\theta_0) &= \sqrt{T} \left[\dot{\Psi}_T(\theta_0) + A_T \right] (\widehat{\theta}_T - \theta_0) \\ \implies -\sqrt{T}\Psi_T(\theta_0) &= \sqrt{T} \underbrace{\left[\dot{\Psi}_T(\theta_0) - \frac{1}{T} \sum_{t=1}^T E_{\theta_0} \left(\frac{\partial}{\partial \theta^T} s(X_t, \theta) \Big|_{\theta_0} \mid \mathcal{F}_{t-1} \right) \right]}_{(*)} \\ &\quad + \frac{1}{T} \sum_{t=1}^T E_{\theta_0} \left(\frac{\partial}{\partial \theta^T} s(X_t, \theta) \Big|_{\theta_0} \mid \mathcal{F}_{t-1} \right) + \underbrace{A_T}_{(**)} (\widehat{\theta}_T - \theta_0) \end{aligned}$$

We will show that the under-braced terms ($K \times K$ matrices) are $o_P(1)$. For the first term $(*)$, note that it involves the average over t of the P_{kjt} terms defined in Condition 3 of the theorem. Clearly, we have that $E_{\theta_0}[P_{kjt} - E_{\theta_0}[P_{kjt} \mid \mathcal{F}_{t-1}] \mid \mathcal{F}_{t-1}] = 0$, and we also have that

$$\begin{aligned} &E_{\theta_0} [|P_{kjt} - E_{\theta_0}[P_{kjt} \mid \mathcal{F}_{t-1}]|] \\ &= E_{\theta_0} \left\{ \left| \frac{\partial}{\partial \theta_j} s_k(X_t, \theta) \Big|_{\theta_0} - E_{\theta_0} \left[\frac{\partial}{\partial \theta_j} s_k(X_t, \theta) \Big|_{\theta_0} \mid \mathcal{F}_{t-1} \right] \right| \right\} \\ &\leq E_{\theta_0} \left\{ \left| \frac{\partial}{\partial \theta_j} s_k(X_t, \theta) \Big|_{\theta_0} \right| \right\} + E_{\theta_0} \left\{ \left| E_{\theta_0} \left[\frac{\partial}{\partial \theta_j} s_k(X_t, \theta) \Big|_{\theta_0} \mid \mathcal{F}_{t-1} \right] \right| \right\} && \text{(Triangle inequality)} \\ &\leq E_{\theta_0} \left\{ \left| \frac{\partial}{\partial \theta_j} s_k(X_t, \theta) \Big|_{\theta_0} \right| \right\} + E_{\theta_0} \left\{ E_{\theta_0} \left[\left| \frac{\partial}{\partial \theta_j} s_k(X_t, \theta) \Big|_{\theta_0} \right| \mid \mathcal{F}_{t-1} \right] \right\} && \text{(Jensen's inequality)} \\ &= 2E_{\theta_0} \left\{ \left| \frac{\partial}{\partial \theta_j} s_k(X_t, \theta) \Big|_{\theta_0} \right| \right\} < \infty && \text{(Condition 3)} \end{aligned}$$

So the assumptions of [Chow \(1965\)](#) (Theorem 5), which is also stated in [Stout \(1974\)](#) (Theorem 3.3.1), are satisfied and we have that

$$\begin{aligned} & \dot{\Psi}_T(\theta_0) - \frac{1}{T} \sum_{t=1}^T E_{\theta_0} \left(\left. \frac{\partial}{\partial \theta^\top} s(X_t, \theta) \right|_{\theta_0} \mid \mathcal{F}_{t-1} \right) \\ &= \frac{1}{T} \sum_{t=1}^T \left[\left. \frac{\partial}{\partial \theta^\top} X_{t,i}(\theta) \right|_{\theta_0} - E_{\theta_0} \left(\left. \frac{\partial}{\partial \theta^\top} s(X_t, \theta) \right|_{\theta_0} \mid \mathcal{F}_{t-1} \right) \right] \xrightarrow{p} 0. \end{aligned}$$

Then for A_T we notice that

$$\left| \left[\ddot{\Psi}_T(\theta^*) \right]_{kji} \right| \leq \frac{1}{T} \sum_{t=1}^T \left| \left. \frac{\partial^2}{\partial \theta_j \partial \theta_i} s_k(X_t, \theta) \right|_{\theta^*} \right| \leq \frac{1}{T} \sum_{t=1}^T \ddot{\psi}(X_t),$$

where the last inequality holds for large T because $\widehat{\theta}_T$ is consistent for θ_0 and the parameter space Θ is an open subset of \mathbb{R}^n which imply that $\widehat{\theta}_T$ is within the neighborhood of θ_0 that satisfies [Condition 4](#) of the theorem with probability that tends to 1, and therefore so will θ^* . Since $\ddot{\psi}(x)$ is integrable, the right hand side above is bounded with probability 1 from the law of large numbers. Then, using Cauchy-Schwarz on $[A_T]_{ki}$ and since $\widehat{\theta}_T$ is consistent for θ_0 , we have that $[A_T]_{ki} \xrightarrow{p} 0$ for all k, i . Therefore, using [Condition 2](#) of the theorem

$$-\sqrt{T}\Psi_T(\theta_0) = \sqrt{T}[V_d + o_P(1)](\widehat{\theta}_T - \theta_0)$$

which, since V_d is invertible, implies asymptotically that

$$\sqrt{T}(\widehat{\theta}_T - \theta_0) \xrightarrow{d} N\left(0, V_d^{-1}V(V_d^{-1})^\top\right).$$

□

[Theorem A.3](#) will be the basis for showing asymptotic normality of our estimators when the propensity score is estimated using a correctly specified parametric propensity score.

Lemma A.1 (Properties of the time series score functions.). *If [Assumption 1](#) holds, and $\psi(w_t, \bar{h}_{t-1}; \gamma)$ are score functions that satisfy [Assumption A.2](#), then*

1. $E_{\gamma_0}[\psi(W_t, \bar{H}_{t-1}; \gamma_0) \mid \mathcal{F}_{t-1}] = 0$, $E_{\gamma_0}[\|\psi(W_t, \bar{H}_{t-1}; \gamma_0)\|] < \infty$, and
2. $E_{\gamma} \left(-\frac{\partial}{\partial \gamma^\top} \psi(W_t, \bar{H}_{t-1}; \gamma) \mid \mathcal{F}_{t-1} \right) = E_{\gamma} \left(\psi(W_t, \bar{H}_{t-1}; \gamma) \psi(W_t, \bar{H}_{t-1}; \gamma)^\top \mid \mathcal{F}_{t-1} \right)$ which in turn implies that $\frac{1}{T} \sum_{t=1}^T E_{\gamma_0} \left(-\frac{\partial}{\partial \gamma^\top} \psi(W_t, \bar{H}_{t-1}; \gamma) \Big|_{\gamma_0} \mid \mathcal{F}_{t-1} \right) \xrightarrow{p} V_{ps}$, for V_{ps} positive definite, symmetric and therefore invertible.

Proof. First, we show that $E_{\gamma_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma_0)\| \right] < \infty$. From Jensen's inequality we have that

$$E_{\gamma_0}^2 \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma_0)\| \right] \leq E_{\gamma_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma_0)\|^2 \right] < \infty,$$

so this part is shown. The remaining of the proof follows steps similar to the ones in the iid setting while conditioning on the corresponding filtration. Since $\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \gamma)$ are the score functions, we have that

$$\begin{aligned} & \boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \gamma) f(W_t = w_t \mid \bar{H}_{t-1} = \bar{h}_{t-1}; \gamma) \\ &= \left[\frac{\partial}{\partial \boldsymbol{\gamma}} \log f(W_t = w_t \mid \bar{H}_{t-1} = \bar{h}_{t-1}; \gamma) \right] f(W_t = w_t \mid \bar{H}_{t-1} = \bar{h}_{t-1}; \gamma) \\ &= \frac{\partial}{\partial \boldsymbol{\gamma}} f(W_t = w_t \mid \bar{H}_{t-1} = \bar{h}_{t-1}; \gamma). \end{aligned} \quad (\text{A.8})$$

Then,

$$\begin{aligned} E_{\gamma}[\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \mid \mathcal{F}_{t-1}] &= E_{\gamma}[\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \mid \bar{H}_{t-M}^*] \\ &= E_{\gamma} \left\{ E_{\gamma}[\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \mid \bar{H}_{t-1}^*] \mid \bar{H}_{t-M}^* \right\} && (\text{Since } \bar{H}_{t-1}^* \supseteq \bar{H}_{t-M}^*) \\ &= E_{\gamma} \left\{ \left[\int \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma) f(W_t = w_t \mid \bar{H}_{t-1}^*) dw_t \right] \mid \bar{H}_{t-M}^* \right\} \\ &= E_{\gamma} \left\{ \left[\int \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma) f(W_t = w_t \mid \bar{H}_{t-1}) dw_t \right] \mid \bar{H}_{t-M}^* \right\} && (\text{Assumption 1}) \\ &= E_{\gamma} \left\{ \left[\int \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma) f(W_t = w_t \mid \bar{H}_{t-1}; \gamma) dw_t \right] \mid \bar{H}_{t-M}^* \right\} \\ &= E_{\gamma} \left\{ \left[\int \frac{\partial}{\partial \boldsymbol{\gamma}^\top} f(W_t = w_t \mid \bar{H}_{t-1}; \gamma) dw_t \right] \mid \bar{H}_{t-M}^* \right\} && (\text{Equation A.8}) \\ &= E_{\gamma} \left\{ \frac{\partial}{\partial \boldsymbol{\gamma}^\top} \left[\int f(W_t = w_t \mid \bar{H}_{t-1}; \gamma) dw_t \right] \mid \bar{H}_{t-M}^* \right\} = 0, \end{aligned}$$

where reversing the integral and derivative is valid using the Leibniz's rule which requires mild regularity conditions (continuity of the propensity score and its partial derivatives with respect to γ). The last equation is equal to zero since the integral of the propensity score over its support is equal to 1, and the derivative to 1 is equal to 0.

To show the second part, we differentiate $E_{\gamma}[\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \mid \mathcal{F}_{t-1}] = 0$ with respect to γ :

$$\begin{aligned} 0 &= \frac{\partial}{\partial \boldsymbol{\gamma}^\top} E_{\gamma}[\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \mid \mathcal{F}_{t-1}] \\ &= \frac{\partial}{\partial \boldsymbol{\gamma}^\top} E_{\gamma} \left\{ E_{\gamma}[\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \mid \bar{H}_{t-1}^*] \mid \bar{H}_{t-M}^* \right\} \\ &= \frac{\partial}{\partial \boldsymbol{\gamma}^\top} E_{\gamma} \left\{ \int \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma) f(W_t = w_t \mid \bar{H}_{t-1}; \gamma) dw_t \mid \bar{H}_{t-M}^* \right\} && (\text{Assumption 1}) \end{aligned}$$

$$\begin{aligned}
&= E_\gamma \left\{ \int \frac{\partial}{\partial \gamma^\top} \left[\boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma) f(W_t = w_t \mid \bar{H}_{t-1}; \gamma) \right] dw_t \mid \bar{H}_{t-M}^* \right\} && \text{(Leibniz's rule)} \\
&= E_\gamma \left\{ \int \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma) \frac{\partial}{\partial \gamma^\top} \left[f(W_t = w_t \mid \bar{H}_{t-1}; \gamma) \right] dw_t \mid \bar{H}_{t-M}^* \right\} \\
&\quad + E_\gamma \left\{ \int \left[\frac{\partial}{\partial \gamma^\top} \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma) \right] f(W_t = w_t \mid \bar{H}_{t-1}; \gamma) dw_t \mid \bar{H}_{t-M}^* \right\} \\
&= E_\gamma \left\{ \int \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma) \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma)^\top f(W_t = w_t \mid \bar{H}_{t-1}; \gamma) dw_t \mid \bar{H}_{t-M}^* \right\} \\
&\hspace{15em} \text{(Equation (A.8))} \\
&\quad + E_\gamma \left\{ \int \left[\frac{\partial}{\partial \gamma^\top} \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma) \right] f(W_t = w_t \mid \bar{H}_{t-1}; \gamma) dw_t \mid \bar{H}_{t-M}^* \right\} \\
&= E_\gamma \left\{ \int \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma) \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma)^\top f(W_t = w_t \mid \bar{H}_{t-1}^*; \gamma) dw_t \mid \bar{H}_{t-M}^* \right\} \\
&\quad + E_\gamma \left\{ \int \left[\frac{\partial}{\partial \gamma^\top} \boldsymbol{\psi}(w_t, \bar{H}_{t-1}; \gamma) \right] f(W_t = w_t \mid \bar{H}_{t-1}^*; \gamma) dw_t \mid \bar{H}_{t-M}^* \right\} \\
&\hspace{15em} \text{(Assumption 1)} \\
&= E_\gamma \left\{ E_\gamma \left[\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma)^\top \mid \bar{H}_{t-1}^* \right] \mid \bar{H}_{t-M}^* \right\} \\
&\quad + E_\gamma \left\{ E_\gamma \left[\frac{\partial}{\partial \gamma^\top} \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \mid \bar{H}_{t-1}^* \right] \mid \bar{H}_{t-M}^* \right\} \\
&= E_\gamma \left[\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma)^\top \mid \bar{H}_{t-M}^* \right] + E_\gamma \left[\frac{\partial}{\partial \gamma^\top} \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \mid \bar{H}_{t-M}^* \right] \\
&\implies E_\gamma \left[\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma)^\top \mid \bar{H}_{t-M}^* \right] = E_\gamma \left[- \frac{\partial}{\partial \gamma^\top} \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma) \mid \bar{H}_{t-M}^* \right].
\end{aligned}$$

From Condition 1b of Assumption A.2 we have the last result. \square

Corollary A.1 (Asymptotic normality of spatio-temporal propensity score parameters). *Consider a propensity score model that satisfies Assumption A.2 and therefore the results of Lemma A.1 hold. Theorem A.3 implies that the MLE of the propensity score parameters are asymptotically normal centered at the true value and with asymptotic variance V_{ps}^{-1} , as in the iid setting.*

Before we state our main theorem we establish a useful Lemma.

Lemma A.2. *Assume that Assumption 1 holds. Let $\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \gamma)$ be the score functions of a propensity score model that satisfies Assumption A.2 as in Lemma A.1 and \mathcal{F}_{t-1} be as above. For*

$$s(\bar{H}_{t-1}, W_t, Y_t; \gamma) = \left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j; \gamma)} \right] N_B(Y_t) - N_{Bt}(F_h^M),$$

it holds that

$$1. E_{\gamma_0} \left[s(\bar{H}_{t-1}, W_t, Y_t; \gamma_0) \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \gamma_0) \mid \mathcal{F}_{t-1} \right] = -E_{\gamma_0} \left[\frac{\partial}{\partial \gamma} s(\bar{H}_{t-1}, W_t, Y_t; \gamma) \Big|_{\gamma_0} \mid \mathcal{F}_{t-1} \right],$$

2. $\frac{\partial}{\partial \gamma_l} s(\bar{h}_{t-1}, w_t, y_t; \gamma) = -N_B(y_t) \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} \right] \sum_{j=t-M+1}^t \psi_l(w_j, \bar{h}_{j-1}; \gamma)$, where we use $\psi_l(w_j, \bar{h}_{j-1}; \gamma)$ to denote the l^{th} element of the $\psi(w_t, \bar{h}_{t-1}; \gamma)$ vector, and
3. similarly $\frac{\partial}{\partial \gamma_m} \frac{\partial}{\partial \gamma_l} s(\bar{h}_{t-1}, w_t, y_t; \gamma)$ is equal to

$$-N_B(y_t) \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} \right] \left\{ \left[\sum_{j=t-M+1}^t \frac{\partial}{\partial \gamma_m} \psi_l(w_j, \bar{h}_{j-1}; \gamma) \right] - \left[\sum_{j=t-M+1}^t \psi_m(w_j, \bar{h}_{j-1}; \gamma) \right] \left[\sum_{j=t-M+1}^t \psi_l(w_j, \bar{h}_{j-1}; \gamma) \right] \right\}$$

Note: $s(\bar{H}_{t-1}, W_t, Y_t; \gamma_0)$ is the term A_{1t} in the proof of Theorem 1.

Proof.

1. We will show it for $M = 1$, and the proof for $M > 1$ is similar. For $M = 1$, $\mathcal{F}_{t-1} = \bar{H}_{t-1}^* = \{\bar{W}_{t-1}, \bar{Y}_T, \bar{X}_T\}$, we consider

$$\begin{aligned} & E_{\gamma_0} [s(\bar{H}_{t-1}, W_t, Y_t; \gamma_0) \psi(W_t, \bar{H}_{t-1}; \gamma_0) \mid \mathcal{F}_{t-1}] \\ &= \int s(\bar{H}_{t-1}, w_t, Y_t; \gamma_0) \psi(w_t, \bar{H}_{t-1}; \gamma_0) f(W_t = w_t \mid \mathcal{F}_{t-1}; \gamma_0) dw_t \\ &= \int s(\bar{H}_{t-1}, w_t, Y_t; \gamma_0) \psi(w_t, \bar{H}_{t-1}; \gamma_0) f(W_t = w_t \mid \bar{H}_{t-1} = \bar{h}_{t-1}; \gamma_0) dw_t \quad (\text{Assumption 1}) \\ &= \int s(\bar{H}_{t-1}, w_t, Y_t; \gamma_0) \frac{\partial}{\partial \gamma} f(W_t = w_t \mid \bar{H}_{t-1} = \bar{h}_{t-1}; \gamma) \Big|_{\gamma_0} dw_t \quad (\text{Equation (A.8)}) \\ &= \int \frac{\partial}{\partial \gamma} [s(\bar{H}_{t-1}, w_t, Y_t; \gamma_0) f(W_t = w_t \mid \bar{H}_{t-1} = \bar{h}_{t-1}; \gamma_0)] dw_t - \\ &\quad - \int \frac{\partial}{\partial \gamma} s(\bar{H}_{t-1}, w_t, Y_t; \gamma) \Big|_{\gamma_0} f(W_t = w_t \mid \bar{H}_{t-1} = \bar{h}_{t-1}; \gamma_0) dw_t \\ &= \int \frac{\partial}{\partial \gamma} [s(\bar{H}_{t-1}, w_t, Y_t; \gamma_0) f(W_t = w_t \mid \mathcal{F}_{t-1}; \gamma_0)] dw_t - \\ &\quad - \int \frac{\partial}{\partial \gamma} s(\bar{H}_{t-1}, w_t, Y_t; \gamma) \Big|_{\gamma_0} f(W_t = w_t \mid \mathcal{F}_{t-1}; \gamma_0) dw_t \quad (\text{Assumption 1}) \\ &= \frac{\partial}{\partial \gamma} E_{\gamma} [s(\bar{H}_{t-1}, W_t, Y_t; \gamma) \mid \mathcal{F}_{t-1}] \Big|_{\gamma_0} - E_{\gamma_0} \left[\frac{\partial}{\partial \gamma} s(\bar{H}_{t-1}, w_t, Y_t; \gamma) \Big|_{\gamma_0} \mid \mathcal{F}_{t-1} \right] \\ &= -E_{\gamma_0} \left[\frac{\partial}{\partial \gamma} s(\bar{H}_{t-1}, w_t, Y_t; \gamma) \Big|_{\gamma_0} \mid \mathcal{F}_{t-1} \right] \end{aligned}$$

where the last equation holds from Equation (A.4). This shows that the expectation is 0, so the derivative is also 0.

Note that at the second line of the proof, we would also need the distribution of Y_t given the filtration \mathcal{F}_{t-1} and the treatment at time period t , $W_t = w_t$. However, given both \mathcal{F}_{t-1} and W_t , the variable Y_t is no longer random, and it is equal to its potential value $Y_t(\overline{W}_{t-1}, w_t)$, where \overline{W}_{t-1} is specified in \mathcal{F}_{t-1} . We refrain from explicitly including this in the proof for simplicity.

2.

$$\begin{aligned}
& \frac{\partial}{\partial \gamma_l} s(\bar{h}_{t-1}, w_t, y_t; \gamma) \\
&= \frac{\partial}{\partial \gamma_l} \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} N_B(y_t) \right] \\
&= N_B(y_t) \left[\prod_{j=t-M+1}^t f_h(w_j) \right] \left[\frac{\partial}{\partial \gamma_l} \frac{1}{\prod_{j=t-M+1}^t e_j(w_j; \gamma)} \right] \\
&= -N_B(y_t) \left[\prod_{j=t-M+1}^t f_h(w_j) \right] \frac{\frac{\partial}{\partial \gamma_l} \prod_{j=t-M+1}^t e_j(w_j; \gamma)}{\left[\prod_{j=t-M+1}^t e_j(w_j; \gamma) \right]^2} \\
&= -N_B(y_t) \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} \right] \frac{\frac{\partial}{\partial \gamma_l} \prod_{j=t-M+1}^t e_j(w_j; \gamma)}{\prod_{j=t-M+1}^t e_j(w_j; \gamma)} \\
&= -N_B(y_t) \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} \right] \sum_{j=t-M+1}^t \frac{\frac{\partial}{\partial \gamma_l} e_j(w_j; \gamma)}{e_j(w_j; \gamma)} \\
&= -N_B(y_t) \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} \right] \sum_{j=t-M+1}^t \psi_{l-1}(w_j, \bar{h}_{j-1}; \gamma). \tag{Equation (A.8)}
\end{aligned}$$

3. Following a similar procedure we have that

$$\begin{aligned}
& \frac{\partial}{\partial \gamma_m} \frac{\partial}{\partial \gamma_l} s(\bar{h}_{t-1}, w_t, y_t; \gamma) \\
&= -N_B(y_t) \left\{ \left[\frac{\partial}{\partial \gamma_m} \prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} \right] \sum_{j=t-M+1}^t \psi_l(w_j, \bar{h}_{j-1}; \gamma) \right. \\
&\quad \left. + \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} \right] \left[\frac{\partial}{\partial \gamma_m} \sum_{j=t-M+1}^t \psi_l(w_j, \bar{h}_{j-1}; \gamma) \right] \right\} \\
&= -N_B(y_t) \left\{ - \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} \right] \left[\sum_{j=t-M+1}^t \psi_m(w_j, \bar{h}_{j-1}; \gamma) \right] \left[\sum_{j=t-M+1}^t \psi_l(w_j, \bar{h}_{j-1}; \gamma) \right] \right. \\
&\quad \left. + \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} \right] \left[\sum_{j=t-M+1}^t \frac{\partial}{\partial \gamma_m} \psi_l(w_j, \bar{h}_{j-1}; \gamma) \right] \right\}
\end{aligned}$$

$$\begin{aligned}
&= -N_B(y_t) \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \gamma)} \right] \left\{ \left[\sum_{j=t-M+1}^t \frac{\partial}{\partial \gamma_m} \psi_l(w_j, \bar{h}_{j-1}; \gamma) \right] \right. \\
&\quad \left. - \left[\sum_{j=t-M+1}^t \psi_m(w_j, \bar{h}_{j-1}; \gamma) \right] \left[\sum_{j=t-M+1}^t \psi_l(w_j, \bar{h}_{j-1}; \gamma) \right] \right\}
\end{aligned}$$

□

Corollary A.2. *Part 1 of Lemma A.1 holds for any function $s(\bar{H}_{t-1}, W_t, Y_t; \gamma)$ for which*

$$E_\gamma [s(\bar{H}_{t-1}, W_t, Y_t; \gamma) \mid \mathcal{F}_{t-1}] = 0.$$

(The proof is identical, hence it is omitted.)

We remind one last result from real analysis which we will use in our theorem. We state it here to avoid unnecessarily complicated notation in the proof of the main theorem. The result extends to multivariate functions.

Remark A.2. *For a function $f : \mathbb{R} \rightarrow \mathbb{R}$ differentiable, if $|f'(x)| \leq \alpha$ for $x \in (x_0 - \epsilon, x_0 + \epsilon)$ and some $\alpha \in \mathbb{R}^+$, then $|f(x)|$ is also bounded on $(x_0 - \epsilon, x_0 + \epsilon)$.*

Proof. The proof is straightforward using Taylor expansion:

$$f(x) = f(x_0) + f'(x^*)(x - x_0) \rightarrow |f(x)| \leq |f(x_0)| + \alpha \epsilon.$$

□

Now we can prove our theorem on asymptotic normality of the causal estimators using propensity scores that are estimated based on a correctly specified propensity score model.

Proof of Theorem 2. We will use Theorem A.3 to show asymptotic normality for the causal estimator based on the estimated propensity score model.

Remember that $\bar{H}_t = \{\bar{W}_t, \bar{Y}_t, \bar{X}_{t+1}\}$. Then $\{\bar{H}_{t-1}, W_t, Y_t\} = \bar{H}_t \setminus \{X_{t-1}\}$ is the set of observed variables until (and including) the t^{th} outcome. Let $\mu \in \mathbb{R}$ and $\gamma \in \mathbb{R}^K$ be the parameters of the propensity score model with score functions $\psi(w_t, \bar{h}_{t-1}; \gamma)$, and define $\theta^\top = (\mu, \gamma^\top)$. Again based on Equation (A.2), we will show the asymptotic normality of the estimator that excludes spatial smoothing. We will then prove that the spatial smoothing does not affect estimation asymptotically because it

converges to zero faster than $T^{-1/2}$. Focusing on the first part of the error, define the $K + 1$ vector

$$s(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}) = \begin{pmatrix} \left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j; \boldsymbol{\gamma})} \right] N_B(Y_t) - N_{Bt}(F_h^M) - \mu \\ \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}) \end{pmatrix} = \begin{pmatrix} A_{1t} - \mu \\ \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}) \end{pmatrix},$$

where A_{1t} is defined in the proof of Theorem 1. We again work with the filtration $\mathcal{F}_t = \bar{H}_{t-M+1}^* = \{\bar{W}_{t-M+1}, \bar{Y}_T, \bar{X}_T\}$. We will show that the conditions of Theorem A.3 hold.

Condition 1a We wish to show the expectation of s conditional on the filtration is 0. Since we showed in the proof of Theorem 1 that

$$\mathbb{E} \left\{ \left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j; \boldsymbol{\gamma})} \right] N_B(Y_t) \mid \mathcal{F}_{t-1} \right\} = N_{Bt}(F_h^M),$$

we have $\boldsymbol{\theta}_0^\top = (\mu_0, \boldsymbol{\gamma}_0^\top) = (0, \boldsymbol{\gamma}_0^\top)$, where $\boldsymbol{\gamma}_0$ represents the true value for the parametric propensity score. Then, based on Lemma A.1, we have that $\mathbb{E}_{\boldsymbol{\theta}_0} [s(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}_0) \mid \mathcal{F}_{t-1}] = 0$. Also, from Jensen's inequality we have that

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\theta}_0}^2 [\|s(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}_0)\|] &\leq \mathbb{E}_{\boldsymbol{\theta}_0} [\|s(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}_0)\|^2] \\ &= \mathbb{E}_{\boldsymbol{\theta}_0}(A_{1t}^2) + \mathbb{E}_{\boldsymbol{\theta}_0} \{ \|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 \} < \infty \end{aligned}$$

where the first term is finite because A_{1t} is bounded as shown in Equation (A.3), and the second term is finite based on Assumption A.2.

Condition 1b Since all terms are under the $\boldsymbol{\theta}_0$ -law, we work with $\mu = \mu_0 = 0$. We have that

$$\begin{aligned} &E_{\boldsymbol{\theta}_0} \left(s(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}_0) s(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}_0)^\top \mid \mathcal{F}_{t-1} \right) \\ &= \begin{bmatrix} E_{\boldsymbol{\theta}_0} [A_{1t}^2 \mid \mathcal{F}_{t-1}] & E_{\boldsymbol{\theta}_0} [A_{1t} \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)^\top \mid \mathcal{F}_{t-1}] \\ E_{\boldsymbol{\theta}_0} [A_{1t} \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0) \mid \mathcal{F}_{t-1}] & E_{\boldsymbol{\theta}_0} [\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0) \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)^\top \mid \mathcal{F}_{t-1}]. \end{bmatrix} \end{aligned}$$

Equation (A.5) implies that $(T - M + 1)^{-1} \sum_{t=M}^T E_{\boldsymbol{\theta}_0} [A_{1t}^2 \mid \mathcal{F}_{t-1}] \xrightarrow{p} v$. In addition, due to Assumption

A.2(1b), we also know that $(T - M + 1)^{-1} \sum_{t=M}^T E_{\boldsymbol{\gamma}_0} \left(\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0) \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)^\top \mid \mathcal{F}_{t-1} \right) \xrightarrow{p} V_{ps}$.

Lastly, Assumption A.3 implies that $(T - M + 1)^{-1} \sum_{t=M}^T E_{\boldsymbol{\theta}_0} [A_{1t} \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0) \mid \mathcal{F}_{t-1}] \xrightarrow{p} u$. Since all the entries of the matrix converge, we are left to show that the resulting matrix is positive definite.

However, since

$$M = \begin{bmatrix} A_{1t}^2 & A_{1t}\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)^\top \\ A_{1t}\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0) & \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)^2 \end{bmatrix}$$

is positive definite (easy to check by taking vector $\boldsymbol{x} \in \mathbb{R}^k$, not all zero, and showing that $\boldsymbol{x}^\top M \boldsymbol{x} > 0$), we have that

$$\begin{bmatrix} v & u^\top \\ u & V_{ps} \end{bmatrix}$$

will also be positive definite.

Condition 1c Take $\epsilon > 0$ and write

$$\begin{aligned} & \frac{1}{T-M+1} \sum_{t=M}^T E_{\theta_0} \left[\|s(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}_0)\|^2 I \left(\|s(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}_0)\| > \epsilon\sqrt{T} \right) \mid \mathcal{F}_{t-1} \right] \\ &= \frac{1}{T-M+1} \sum_{t=M}^T E_{\theta_0} \left[(A_{1t}^2 + \|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2) I \left(A_{1t}^2 + \|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T \right) \mid \mathcal{F}_{t-1} \right] \\ &= \frac{1}{T-M+1} \sum_{t=M}^T E_{\theta_0} \left[A_{1t}^2 I \left(A_{1t}^2 + \|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T \right) \mid \mathcal{F}_{t-1} \right] \\ & \quad + \frac{1}{T-M+1} \sum_{t=M}^T E_{\theta_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T - A_{1t}^2 \right) \mid \mathcal{F}_{t-1} \right] \end{aligned}$$

We start with the second term: Since A_{1t}^2 cannot exceed $(\delta_W^M \delta_Y + \delta_Y)^2$ based on Equation (A.3), we have that

$$\begin{aligned} & \frac{1}{T-M+1} \sum_{t=M}^T E_{\theta_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T - A_{1t}^2 \right) \mid \mathcal{F}_{t-1} \right] \\ & \leq \frac{1}{T-M+1} \sum_{t=M}^T E_{\theta_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\| > \sqrt{\epsilon^2 T - (\delta_W^M \delta_Y + \delta_Y)^2} \right) \mid \mathcal{F}_{t-1} \right] \xrightarrow{p} 0, \end{aligned}$$

based on Assumption A.2 and since $\delta_W^M \delta_Y + \delta_Y$ is fixed.

For the first term, since $I(A_{1t}^2 + \|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2) > \epsilon^2 T$ implies that at least one of A_{1t}^2 and $\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2$ is greater than $\epsilon^2 T/2$, we have that

$$I(A_{1t}^2 + \|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2) > \epsilon^2 T \leq I(A_{1t}^2 > \epsilon^2 T/2) + I(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2).$$

This leads to

$$E_{\theta_0} \left[A_{1t}^2 I \left(A_{1t}^2 + \|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T \right) \mid \mathcal{F}_{t-1} \right]$$

$$\leq E_{\theta_0} \left[A_{1t}^2 I \left(A_{1t}^2 > \epsilon^2 T/2 \right) \mid \mathcal{F}_{t-1} \right] + E_{\theta_0} \left[A_{1t}^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2 \right) \mid \mathcal{F}_{t-1} \right].$$

In the proof of Theorem 1 we have already shown that because A_{1t} is bounded we have that

$$\frac{1}{T-M+1} \sum_{t=M}^T E_{\theta_0} \left[A_{1t}^2 I \left(|A_{1t}| > \frac{\epsilon}{\sqrt{2}} \sqrt{T} \right) \mid \mathcal{F}_{t-1} \right] \xrightarrow{p} 0,$$

and we want to show that

$$\frac{1}{T-M+1} \sum_{t=M}^T E_{\theta_0} \left[A_{1t}^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2 \right) \mid \mathcal{F}_{t-1} \right] \xrightarrow{p} 0.$$

We write

$$\begin{aligned} & E_{\theta_0} \left[A_{1t}^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2 \right) \mid \mathcal{F}_{t-1} \right] \\ &= E_{\theta_0} \left[A_{1t}^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2 \right) \mid A_{1t}^2 \leq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2, \mathcal{F}_{t-1} \right] \times \\ &\quad \times P(A_{1t}^2 \leq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2 \mid \mathcal{F}_{t-1}) + \\ &\quad + E_{\theta_0} \left[A_{1t}^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2 \right) \mid A_{1t}^2 \geq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2, \mathcal{F}_{t-1} \right] \times \\ &\quad \times P(A_{1t}^2 \geq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2 \mid \mathcal{F}_{t-1}) \\ &\leq E_{\theta_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2 \right) \mid A_{1t}^2 \leq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2, \mathcal{F}_{t-1} \right] \times \\ &\quad \times P(A_{1t}^2 \leq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2 \mid \mathcal{F}_{t-1}) + \\ &\quad + E_{\theta_0} \left[A_{1t}^2 I \left(A_{1t}^2 > \epsilon^2 T/2 \right) \mid A_{1t}^2 \geq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2, \mathcal{F}_{t-1} \right], \end{aligned}$$

where again the average over time of the last term will be converging to zero in probability since A_{1t} is bounded and using similar arguments. Using the law of total expectation we can write

$$\begin{aligned} & E_{\theta_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2 \right) \mid \mathcal{F}_{t-1} \right] \\ &= E_{\theta_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2 \right) \mid A_{1t}^2 \leq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2, \mathcal{F}_{t-1} \right] \times \\ &\quad \times P(A_{1t}^2 \leq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2 \mid \mathcal{F}_{t-1}) + \\ &\quad E_{\theta_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2 \right) \mid A_{1t}^2 \geq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2, \mathcal{F}_{t-1} \right] \times \\ &\quad \times P(A_{1t}^2 \geq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2 \mid \mathcal{F}_{t-1}) \\ &\geq E_{\theta_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 I \left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2 \right) \mid A_{1t}^2 \leq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2, \mathcal{F}_{t-1} \right] \times \\ &\quad \times P(A_{1t}^2 \leq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2 \mid \mathcal{F}_{t-1}). \end{aligned}$$

Since all the terms in the expectations are positive and since (from Assumption A.2) we have that

$$\frac{1}{T-M+1} \sum_{t=M}^T E_{\theta_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 I\left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2\right) \mid \mathcal{F}_{t-1} \right] \xrightarrow{p} 0$$

we also have that

$$\begin{aligned} & \frac{1}{T-M+1} \sum_{t=M}^T \left\{ \right. \\ & E_{\theta_0} \left[\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 I\left(\|\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0)\|^2 > \epsilon^2 T/2\right) \mid A_{1t}^2 \leq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2, \mathcal{F}_{t-1} \right] \times \\ & \left. \times P(A_{1t}^2 \leq \|\boldsymbol{\psi}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})\|^2 \mid \mathcal{F}_{t-1}) \right\} \xrightarrow{p} 0 \end{aligned}$$

which completes the proof that Condition 1c holds.

Condition 2 We denote $\boldsymbol{\theta}^\top = (\theta_1, \theta_2, \dots, \theta_{K+1}) = (\mu, \boldsymbol{\gamma}^\top)$ and use $s_k(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta})$ to denote the k^{th} entry of the $s(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta})$ vector. We note that

$$\frac{\partial}{\partial \boldsymbol{\theta}^\top} s(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}) = \begin{bmatrix} -1 & \frac{\partial}{\partial \boldsymbol{\gamma}^\top} s_1(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}) \\ \mathbf{0} & \frac{\partial}{\partial \boldsymbol{\gamma}^\top} \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}) \end{bmatrix}$$

Lemma A.1 implies that $(T-M+1)^{-1} \sum_{t=M}^T E_{\boldsymbol{\gamma}_0} \left[\frac{\partial}{\partial \boldsymbol{\gamma}^\top} \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}) \Big|_{\boldsymbol{\gamma}_0} \mid \mathcal{F}_{t-1} \right] \rightarrow -V_{ps}$ (invertible).

Assumption A.3 and Lemma A.2 imply that

$$\begin{aligned} & (T-M+1)^{-1} \sum_{t=M}^T E_{\theta_0} \left(\frac{\partial}{\partial \boldsymbol{\gamma}^\top} s_1(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}) \Big|_{\theta_0} \mid \mathcal{F}_{t-1} \right) \\ & = -(T-M+1)^{-1} \sum_{t=M}^T E_{\theta_0} \left[s_1(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}_0) \boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma}_0) \mid \mathcal{F}_{t-1} \right] \\ & \xrightarrow{p} -u^\top. \end{aligned}$$

Putting these together we have that

$$(T-M+1)^{-1} \sum_{t=M}^T E_{\theta_0} \left[\frac{\partial}{\partial \boldsymbol{\theta}^\top} s(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}) \Big|_{\theta_0} \mid \mathcal{F}_{t-1} \right] \xrightarrow{p} \begin{bmatrix} -1 & -u \\ 0 & -V_{ps} \end{bmatrix}.$$

Since V_{ps} is invertible and the first row is the only one to have a non-zero first element we have that this limit matrix is invertible.

Condition 3 We want to show that for all $k, j = 1, 2, \dots, K + 1$, if we use P_{kjt} to denote

$$P_{kjt} = \frac{\partial}{\partial \theta_j} s_k(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}_0},$$

then $E_{\boldsymbol{\theta}_0} |P_{kjt}| < \infty$, and there exists $0 < r_{kj} \leq 2$ such that

$$\sum_{t=M}^T \frac{1}{t^{r_{kj}}} E_{\boldsymbol{\theta}_0} [|P_{kjt} - E_{\boldsymbol{\theta}_0}(P_{kjt} | \mathcal{F}_{t-1})|^{r_{kj}} | \mathcal{F}_{t-1}] \xrightarrow{P} 0.$$

For $k, j \geq 2$, this is given by Condition 2 of Assumption A.2. For $j = 1$ and $k \geq 2$, we have that

$$\frac{\partial}{\partial \theta_1} s_k(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}) = 0,$$

so the result holds for any r_{k1} . Similarly, for $k = j = 1$, we have that

$$\frac{\partial}{\partial \theta_1} s_1(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}) = -1,$$

so the result holds for a value $r_{11} \in (1, 2]$. Therefore, it is left to show that it holds for $k = 1$ and $j \geq 2$.

For $k = 1$ and $j \geq 2$, the condition that there exists $0 < r_{1j} \leq 2$ such that

$$\sum_{t=M}^T \frac{1}{t^{r_{1j}}} E_{\boldsymbol{\theta}_0} [|P_{1jt} - E_{\boldsymbol{\theta}_0}(P_{1jt} | \mathcal{F}_{t-1})|^{r_{1j}} | \mathcal{F}_{t-1}] \xrightarrow{P} 0$$

is given by Assumption A.3. So we are left to show that $E(|P_{1jt}|) < \infty$. Lemma A.2 implies that

$$\begin{aligned} P_{1jt} &= \frac{\partial}{\partial \gamma_{j-1}} s_1(\bar{H}_{t-1}, W_t, Y_t; \boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}_0} \\ &= -N_B(Y_t) \left[\prod_{t'=t-M+1}^t \frac{f_h(W_{t'})}{e_{t'}(W_{t'}; \boldsymbol{\gamma}_0)} \right] \sum_{t'=t-M+1}^t \boldsymbol{\psi}_{j-1}(W_{t'}, \bar{H}_{t'-1}; \boldsymbol{\gamma}_0) \\ \implies |P_{1jt}| &\leq \delta_Y \delta_W^M \sum_{t'=t-M+1}^t \left| \boldsymbol{\psi}_{j-1}(W_{t'}, \bar{H}_{t'-1}; \boldsymbol{\gamma}_0) \right| \\ \implies E_{\boldsymbol{\theta}_0} |P_{1jt}| &\leq \delta_Y \delta_W^M \sum_{t'=t-M+1}^t E_{\boldsymbol{\gamma}_0} \left| \boldsymbol{\psi}_{j-1}(W_{t'}, \bar{H}_{t'-1}; \boldsymbol{\gamma}_0) \right|. \end{aligned}$$

Since

$$E_{\boldsymbol{\gamma}_0}^2 \left| \boldsymbol{\psi}_{j-1}(W_{t'}, \bar{H}_{t'-1}; \boldsymbol{\gamma}_0) \right| \leq E_{\boldsymbol{\gamma}_0} \left[\boldsymbol{\psi}_{j-1}(W_{t'}, \bar{H}_{t'-1}; \boldsymbol{\gamma}_0)^2 \right] \quad (\text{Jensen's inequality})$$

$$\leq E_{\gamma_0} [\|\boldsymbol{\psi}_{j-1}(W_{t'}, \bar{H}_{t'-1}; \boldsymbol{\gamma}_0)\|^2] < \infty, \quad (\text{Assumption A.2})$$

we have that $E|P_{1jt}| < \infty$.

Condition 4 We want to show that there exists integrable function $\ddot{\psi}(x)$ which dominates the second partial derivatives of $s(\bar{h}_{t-1}, w_t, y_t; \boldsymbol{\theta})$ in a neighborhood of $\boldsymbol{\theta}_0$ for all $(w_t, \bar{h}_{t-1}, y_t)$. We consider derivatives of $s(\bar{h}_{t-1}, w_t, y_t; \boldsymbol{\theta})$ with respect to θ_m, θ_l . For $k, m, l \geq 2$,

$$\frac{\partial}{\theta_m} \frac{\partial}{\theta_l} s_k(\bar{h}_{t-1}, w_t, y_t; \boldsymbol{\theta}) = \frac{\partial}{\gamma_{m-1}} \frac{\partial}{\gamma_{l-1}} \boldsymbol{\psi}_{k-1}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})$$

where $\boldsymbol{\psi}_{k-1}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma})$ is the $k-1$ entry of the $\boldsymbol{\psi}(W_t, \bar{H}_{t-1}; \boldsymbol{\gamma})$ vector. From Condition 3 of Assumption A.2, we know that the above is dominated by an integrable function. For $k \geq 2$ and if $l = 1$ or $m = 1$ we have that the second partial derivative is equal to 0, since

$$\frac{\partial}{\partial \theta_1} s_k(\bar{h}_{t-1}, w_t, y_t; \boldsymbol{\theta}) = \frac{\partial}{\partial \mu} \boldsymbol{\psi}_{k-1}(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma}) = 0.$$

So for $k \geq 2$, all second partial derivatives are dominated by the function in Condition 3 of Assumption A.2. Then, for $k = 1$, if at least one of $l = 1$ or $m = 1$ we have that the second partial derivative is also zero, since

$$\frac{\partial}{\partial \theta_1} s_1(\bar{h}_{t-1}, w_t, y_t; \boldsymbol{\theta}) = \frac{\partial}{\partial \mu} s_1(\bar{h}_{t-1}, w_t, y_t; \boldsymbol{\theta}) = -1.$$

So we need to show it only for $k = 1$, and $l, m \geq 2$. From Lemma A.2 we have that

$$\begin{aligned} & \frac{\partial^2}{\partial \theta_m \partial \theta_l} s_1(\bar{h}_{t-1}, w_t, y_t; \boldsymbol{\theta}) \\ &= -N_B(y_t) \left[\prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \boldsymbol{\gamma})} \right] \left\{ \left[\sum_{j=t-M+1}^t \frac{\partial}{\partial \gamma_{m-1}} \boldsymbol{\psi}_{l-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}) \right] \right. \\ & \quad \left. - \left[\sum_{j=t-M+1}^t \boldsymbol{\psi}_{m-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}) \right] \left[\sum_{j=t-M+1}^t \boldsymbol{\psi}_{l-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}) \right] \right\} \end{aligned}$$

Because of Assumption A.1(a) and Assumption 2 we have that

$$|N_B(y_t)| \leq \delta_Y \quad \text{and} \quad 0 \leq \prod_{j=t-M+1}^t \frac{f_h(w_j)}{e_j(w_j; \boldsymbol{\gamma})} \leq \delta_W^M,$$

which implies that

$$\begin{aligned}
\left| \frac{\partial^2}{\partial \theta_m \partial \theta_l} s_1(\bar{h}_{t-1}, w_t, y_t; \boldsymbol{\theta}) \right| &\leq \delta_Y \delta_W^M \left| \sum_{j=t-M+1}^t \frac{\partial}{\partial \gamma_{m-1}} \psi_{l-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}) \right| + \\
&\quad + \delta_Y \delta_W^M \left| \sum_{j,j'=t-M+1}^t \psi_{m-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}) \psi_{l-1}(w_{j'}, \bar{h}_{j'-1}; \boldsymbol{\gamma}) \right| \\
&\leq \sum_{j=t-M+1}^t \delta_Y \delta_W^M \left| \frac{\partial}{\partial \gamma_{m-1}} \psi_{l-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}) \right| + \\
&\quad + \sum_{j,j'=t-M+1}^t \delta_Y \delta_W^M \left| \psi_{m-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}) \psi_{l-1}(w_{j'}, \bar{h}_{j'-1}; \boldsymbol{\gamma}) \right|
\end{aligned}$$

We work first with the first term. Since the summation is over M terms with M finite, we only need to study the quantity in the absolute value. We know from Assumption A.2 that the second partial derivatives of $\psi(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})$ are dominated by $\ddot{\psi}(w_t, \bar{h}_{t-1})$ in a neighborhood of $\boldsymbol{\gamma}_0$. Assume that this neighborhood is the ϵ -ball around $\boldsymbol{\gamma}_0$ (this always exists since a neighborhood is an open set around $\boldsymbol{\gamma}_0$). Then, from Remark A.2 we know that

$$\left| \frac{\partial}{\partial \gamma_{m-1}} \psi_{l-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}) \right| \leq \left| \frac{\partial}{\partial \gamma_{m-1}} \psi_{l-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}) \right|_{\boldsymbol{\gamma}_0} + \epsilon K \ddot{\psi}(w_t, \bar{h}_{t-1}),$$

where the K appears because we consider all K second partial derivatives which are all bounded by $\ddot{\psi}$. From Assumption A.2(2), we have that the quantity on the right has finite expectation and is fixed in $\boldsymbol{\gamma}$. Therefore, it is an integrable function that dominates the first partial derivatives of $\psi(w_t, \bar{h}_{t-1}; \boldsymbol{\gamma})$ in a neighborhood of $\boldsymbol{\gamma}_0$ for all l, m . Denote the maximum of these functions over l, m by $\ddot{\psi}_1(w_t, \bar{h}_{t-1})$.

We now turn our attention to the second term. Since (using again Remark A.2)

$$\left| \psi_{m-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}) \right| \leq \left| \psi_{m-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}_0) \right| + \epsilon K \ddot{\psi}_1(w_t, \bar{h}_{t-1}),$$

and $E_{\boldsymbol{\gamma}_0} [|\psi_{m-1}(w_j, \bar{h}_{j-1}; \boldsymbol{\gamma}_0)|] < \infty$ from Assumption A.2, we have that this quantity is also dominated by an integrable function that is constant in $\boldsymbol{\gamma}$. Denote the maximum of these functions over m as $\ddot{\psi}_2(w_t, \bar{h}_{t-1})$.

Putting these together we have that

$$\left| \frac{\partial^2}{\partial \theta_m \partial \theta_l} s_1(\bar{h}_{t-1}, w_t, y_t; \boldsymbol{\theta}) \right| \leq M \delta_Y \delta_W^M \ddot{\psi}_1(w_t, \bar{h}_{t-1}) + M^2 \delta_Y \delta_W^M \left[\ddot{\psi}_2(w_t, \bar{h}_{t-1}) \right]^2,$$

where the right hand side is integrable. By defining taking the maximum of the quantity on the right

hand side and $\ddot{\psi}(w_t, \bar{h}_{t-1})$ for each (w_t, \bar{h}_{t-1}) we have that the condition holds using this new integrable function.

Consistency of the solution The last condition of Theorem A.3 that we need to show is that the solution to $\sum_{t=M}^T s(\bar{H}_{t-1}, W_t, Y_t; \theta) = 0$ is consistent for θ_0 . Since the estimator of the propensity score parameters based on the score functions are consistent, we only need to show that the solution to $\sum_{t=M}^T s_1(\bar{H}_{t-1}, W_t, Y_t; \theta) = 0$ is consistent for $\mu_0 = 0$.

Since the estimator based on the true propensity score was shown to be consistent in Theorem 1, the propensity score estimators $\hat{\gamma}$ are consistent for γ , $s_1(\bar{H}_{t-1}, W_t, Y_t; \theta)$ is a continuous function of the propensity score which is itself continuous in γ , using Slutsky's theorem we have that the solution to $\sum_{t=M}^T s_1(\bar{H}_{t-1}, W_t, Y_t; \theta) = 0$ using the estimated propensity score parameters is also consistent.

Asymptotic normality of the estimator without spatial smoothing Since the conditions of Theorem A.3 are satisfied, we have that the solution $\hat{\theta}_T$ to $\sum_{t=M}^T s(\bar{H}_{t-1}, W_t, Y_t; \theta) = 0$ are asymptotically normal with

$$\sqrt{T} \left(\hat{\theta}_T - \theta_0 \right) \xrightarrow{d} N(0, V_\theta),$$

where $V_\theta = A^{-1}B(A^{-1})^T$ for

$$A = \begin{bmatrix} -1 & -u^T \\ \mathbf{0}_K & -V_{ps} \end{bmatrix} \quad \text{and} \quad B = \begin{bmatrix} v & u^T \\ u & V_{ps} \end{bmatrix}. \quad (\text{A.9})$$

As a result, focusing on the first entry of $\hat{\theta}$ and since $\mu_0 = 0$, we have that

$$\underbrace{\sqrt{T} \left\{ (T - M + 1)^{-1} \sum_{t=M}^T \left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j; \gamma)} \right] N_B(Y_t) - N_B(F_h^M) \right\}}_{\text{estimator without spatial smoothing}} \rightarrow N(0, v^e),$$

where $v^e = [V_\theta]_{11}$ is the (1, 1) entry of V_θ .

Asymptotic normality of the estimator with spatial smoothing To prove the asymptotic normality of the estimator with spatial smoothing (our proposed estimator in Equation (8)), we again decompose the estimation error in two components like in Equation (A.2) for the proof of Theorem 1. We write

$$err_t = \underbrace{\left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j; \gamma)} \right] N_B(Y_t) - N_{Bt}(F_h^M)}_{A_{1t}}$$

$$+ \underbrace{\left[\prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j; \gamma)} \right] \left[\int_B \sum_{s \in S_{Y_t}} K_{b_T}(\omega, s) d\omega - N_B(Y_t) \right]}_{A_{2t}},$$

where we use the parametric propensity score. We showed the asymptotic normality based on A_{1t} , so we are left to show that $\sqrt{T} \left((T - M + 1)^{-1} \sum_{t=M}^T A_{2t} \right) \xrightarrow{p} 0$. In the proof of Theorem 1 we already showed that the above result holds. The proof there can be directly used here also if the known propensity score is used (instead of the estimated one). By re-defining the terms c_t defined there to use the estimated propensity score as

$$c_t = \prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j; \gamma)},$$

it suffices to show that c_t is bounded, and the steps of the proof with the known propensity score will follow identically. But since the propensity score $e_t(w; \gamma)$ is continuous in γ (since it is differentiable), the function $1/x$ is continuous for $x > 0$, and $f_h(w_j)/e_j(w_j; \gamma_0) \leq \delta_W$ then c_t will be bounded in a neighborhood of γ_0 . And since $\hat{\gamma} \xrightarrow{p} \gamma_0$, $\hat{\gamma}$ will be in the neighborhood of γ_0 with probability 1 as T increases, so c_t will be bounded.

Putting these results together we have asymptotic normality of the spatially smoothed estimator and

$$\sqrt{T} \left(\hat{N}_B(F_h^M) - N_B(F_h^M) \right) \xrightarrow{d} N(0, v^e).$$

□

Proof of Theorem 3. The asymptotic variance v^e corresponds to the (1, 1) entry of the matrix $A^{-1}B(A^{-1})^\top$, where A, B are defined in Equation (A.9).

$$\begin{aligned} A^{-1}B(A^{-1})^\top &= \begin{bmatrix} 1 & u^\top \\ \mathbf{0}_K & V_{ps} \end{bmatrix}^{-1} \begin{bmatrix} v & u^\top \\ u & V_{ps} \end{bmatrix} \left\{ \begin{bmatrix} 1 & u^\top \\ \mathbf{0}_K & V_{ps} \end{bmatrix}^{-1} \right\}^\top \\ &= \begin{bmatrix} 1 & -u^\top V_{ps}^{-1} \\ \mathbf{0}_K & V_{ps}^{-1} \end{bmatrix} \begin{bmatrix} v & u^\top \\ u & V_{ps} \end{bmatrix} \begin{bmatrix} 1 & \mathbf{0}_K \\ -V_{ps}^{-1}u & V_{ps}^{-1} \end{bmatrix} \\ &= \begin{bmatrix} v - u^\top V_{ps}^{-1}u & \mathbf{0}_K^\top \\ \dots & \dots \end{bmatrix} \begin{bmatrix} 1 & \dots \\ -V_{ps}^{-1}u & \dots \end{bmatrix} \\ &= \begin{bmatrix} v - u^\top V_{ps}^{-1}u & \dots \\ \dots & \dots \end{bmatrix} \end{aligned}$$

so $v^e = v - u^\top V_{ps}^{-1}u$, and since V_{ps} is positive definite and therefore V_{ps}^{-1} is positive definite we have that $u^\top V_{ps}^{-1}u \geq 0$ and $v^e \leq v$. □

B.5 Asymptotics for an increasing number of independent regions

All the asymptotic results that have been discussed up to now correspond to the scenario where 1 region is observed repeatedly over time, and the asymptotic properties are derived when the number of time periods T increases to infinity. However, there might also be interest in situations where the number of time periods is fixed, but there exist an increasing number of independent-acting regions.

Here we consider this related but separate scenario. We start by defining relevant estimands in this setting, ensuring that these new estimands are as closely comparable to the estimands in the manuscript. We propose similar estimators, and derive the asymptotic properties of the new estimators when the number of regions R goes to infinity.

B.5.1 Estimands for independently-acting regions For this scenario, we decompose the treatments, potential outcomes, outcomes, and history over all the regions to region-specific components and write $w_t = (w_{1t}, w_{2t}, \dots, w_{Rt})$, $\bar{w}_t = (\bar{w}_{1t}, \bar{w}_{2t}, \dots, \bar{w}_{Rt})$, $Y_t(\bar{w}_t) = (Y_{1t}(\bar{w}_t), Y_{2t}(\bar{w}_t), \dots, Y_{Rt}(\bar{w}_t))$, $Y_t = (Y_{1t}, Y_{2t}, \dots, Y_{Rt})$, and $\bar{H}_t = (\bar{H}_{1t}, \bar{H}_{2t}, \dots, \bar{H}_{Rt})$, where $\bar{H}_{rt} = \{\bar{W}_{rt}, \bar{Y}_{rt}, \bar{X}_{r(t+1)}\}$. We make the following assumption that describes that the regions do not interfere spatially, and that treatment assignment is local within regions:

Assumption A.4 (Independently acting spatial regions). *We assume the following:*

1. For \bar{w}_t, \bar{w}'_t such that $\bar{w}_{rt} = \bar{w}'_{rt}$, we have that $Y_{rt}(\bar{w}_t) = Y_{rt}(\bar{w}'_t)$ (and a similar assumption for the time-varying covariates), and
2. the treatment assignment of region r at time t does not depend on unobserved potential outcomes or potential time varying covariates, nor on any information from other regions, denoted as $W_{rt} \perp \perp \bar{H}_{t-1}, \bar{Y}_T, \bar{X}_T \mid \bar{H}_{r(t-1)}$.

This assumption allows us to denote potential outcomes using their region-specific treatments *only*, and write $Y_t(\bar{w}_t) = (Y_{1t}(\bar{w}_{1t}), Y_{2t}(\bar{w}_{2t}), \dots, Y_{Rt}(\bar{w}_{Rt}))$. It also allows us to think of the R regions as completely separately acting regions, as outcomes, covariates and treatments of one region do not depend on any information of any other region. Based on this assumption, we can use \bar{Y}_{rT} to denote the collection of potential outcomes for region r over all time periods and for any regional treatment path (and similarly for covariates).

For the purpose of this section *only*, we also assume the temporal carryover effect is limited to up to some lag M_Y . Specifically, we assume that the outcome at time t can only depend on treatments during the preceding M_Y time periods, formalized as

Assumption A.5 (Limited temporal carryover effect). *There exists positive integer M_Y such that for $\bar{w}_{rt}, \bar{w}'_{rt}$ for which $w_{r\tau} = w'_{r\tau}$ for all $\tau = t - M_Y + 1, \dots, t - 1, t$, it holds that $Y_{rt}(\bar{w}_{rt}) = Y_{rt}(\bar{w}'_{rt})$.*

We start by defining region and time specific estimands that are as closely related to the estimands defined in Section 3. We again focus on point pattern treatments and outcomes and on estimands that represent the number of outcome active locations in each region. For simplicity we focus on the scenario where the temporal carryover lag M_Y and the intervention length M are both equal to 1, but we note that the results would also follow in all scenarios where $M \geq M_Y$.

Let F_h be a stochastic treatment assignment that is constant across regions. The stochastic intervention can depend on baseline covariates of the regions, but we refrain for explicitly denoting that for simplicity. We define the expected number of outcome active locations at region r at time t as

$$N_{rt}(F_h) = \int_{w_{rt}} N_r(Y_{rt}(\overline{W}_{r(t-1)}, w_{rt})) dF_h(w_{rt}) = \int_{w_{rt}} N_r(Y_{rt}(w_{rt})) dF_h(w_{rt}),$$

where we define the estimand as in the first equation to be more closely related to the estimands in §3, and the second equation holds because of Assumption A.5 for $M_Y = 1$. We specify region-specific estimands, averaged over time, as

$$N_r(F_h) = \frac{1}{T} \sum_{t=1}^T N_{rt}(F_h),$$

and estimands averaged over region and time as

$$N(F_h) = \frac{1}{R} \sum_{r=1}^R N_r(F_h) = \frac{1}{R} \sum_{r=1}^R \frac{1}{T} \sum_{t=1}^T N_{rt}(F_h).$$

B.5.2 Estimators for independently-acting regions Like in Section 4, assume that F_h admits density f_h . Based on Assumption A.4, we can separate the treatment assignment over all regions to the treatment assignment of each region separately, as

$$e_t(w_t) = \prod_{r=1}^R e_{rt}(w_{rt})$$

where $e_{rt}(w_{rt}) = f(W_{rt} = w_{rt} \mid \overline{H}_{r(t-1)})$ is the region-specific propensity score. We propose corresponding region and time-specific estimator

$$\widehat{N}_{rt}(F_h) = \frac{f_h(W_{rt})}{e_{rt}(W_{rt})} N_r(Y_{rt}),$$

where we use $N_r(Y_{rt})$ to denote the number of outcome active locations in the observed outcome for region r at time t . We also propose the corresponding estimators averaged over time and over regions as

$$\widehat{N}_r(F_h) = \frac{1}{T} \sum_{t=1}^T \widehat{N}_{rt}(F_h) \quad \text{and} \quad \widehat{N}(F_h) = \frac{1}{R} \sum_{r=1}^R \widehat{N}_r(F_h).$$

B.5.3 Consistency and asymptotic normality for independent-acting regions We will show the consistency and asymptotic normality of these estimators when the propensity score is known for an increasing number of independently acting regions. The proof here follows closely the proof in [Papadogeorgou et al. \(2019\)](#) for weighting estimators under a known propensity score and for stochastic interventions. We do not show the asymptotic properties of an estimator based on a correctly specified parametric propensity score, since, once the baseline conditions for the known propensity score are established, the proof for the estimated propensity score would resemble the corresponding proof in [Papadogeorgou et al. \(2019\)](#).

To establish the asymptotic properties for an increasing number of regions, we first assume that our observed regions are a random sample from some super-population of regions. Let $(\overline{Y}_{rT}, \overline{X}_{rT}, \overline{W}_{rT})$ be a draw from a super-population distribution F^{sp} . We assume that Assumption A.4 holds over F^{sp} and we make the following super-population positivity assumption for the independent regions (which resembles the one in Assumption 2):

Assumption A.6 (Positivity of treatment assignment in the super-population). *There exists δ_W such that $e_{rt}(w_{rt}) > \delta_W \cdot f_h(w_{rt})$ for all treatment point patterns w_{rt} .*

We also assume that there is a bounded number of outcome active locations within each region, similarly to Assumption A.1(a), but region-specific:

Assumption A.7. *There exists $\delta_Y > 0$ such that $N_r(Y_{rt}(\overline{w}_{rt})) < \delta_Y$ with probability 1 over F^{sp} , where \overline{w}_{rt} is any possible treatment path.*

Theorem A.4. *If Assumptions A.4, A.5 and A.6 hold, then, for $R \rightarrow \infty$, $\widehat{N}(F_h)$ is consistent for $N(F_h)$ and $\sqrt{R} \left(\widehat{N}(F_h) - N(F_h) \right) \rightarrow N(0, \sigma^2)$, for some $\sigma^2 > 0$, where $N(F_h)$ is the super-population estimand defined as $N(F_h) = E_{F^{sp}} [N_r(F_h)]$.*

Proof. Let $\mathcal{D}_r = (\overline{W}_T, \overline{Y}_T, \overline{X}_T)$ denote all observed data for region r , and $\mathcal{D} = (\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_R)$. We define

$$\psi_r(\mathcal{D}_r; \mu) = \left(\frac{1}{T} \sum_{t=1}^T \frac{f_h(W_{rt})}{e_{rt}(W_{rt})} N_r(Y_{rt}) \right) - \mu$$

and $\Psi_R(\mathcal{D}; \mu) = \sum_{r=1}^T \psi_r(\mathcal{D}_r; \mu)$. Then obviously the estimator $\hat{\mu} = \hat{N}(F_h)$ is the solution to $\Psi_R(\mathcal{D}; \mu) = 0$. Then we calculate the solution to $\Psi^{sp}(\mu) = E_{F^{sp}}(\psi_r(\mathcal{D}_r; \mu)) = 0$ which is equal to

$$\begin{aligned}
\mu_0 &= E_{F^{sp}} \left[\frac{1}{T} \sum_{t=1}^T \frac{f_h(W_{rt})}{e_{rt}(W_{rt})} N_r(Y_{rt}) \right] \\
&= \frac{1}{T} \sum_{t=1}^T E_{F^{sp}} \left[\frac{f_h(W_{rt})}{e_{rt}(W_{rt})} N_r(Y_{rt}) \right] \\
&= \frac{1}{T} \sum_{t=1}^T \int_{\bar{\mathcal{Y}}_{rT}, \bar{\mathcal{X}}_{rT}} \int_{w_{r1}} \int_{w_{r2}} \cdots \int_{w_{rt}} \frac{f_h(w_{rt})}{e_{rt}(w_{rt})} N_r(Y_{rt}(\bar{\mathbf{w}}_{rt})) dF^{sp}(\bar{\mathcal{Y}}_{rT}, \bar{\mathcal{X}}_{rT}, \bar{\mathbf{w}}_{rt}) \\
&= \frac{1}{T} \sum_{t=1}^T \int_{\bar{\mathcal{Y}}_{rT}, \bar{\mathcal{X}}_{rT}} \int_{w_{r1}} \int_{w_{r2}} \cdots \int_{w_{r(t-1)}} \left[\int_{w_{rt}} \frac{f_h(w_{rt})}{e_{rt}(w_{rt})} N_r(Y_{rt}(\bar{\mathbf{w}}_{rt})) f_{W_{rt}}(w_{rt} | \bar{\mathbf{W}}_{r(t-1)} = \bar{\mathbf{w}}_{r(t-1)}, \bar{\mathcal{Y}}_{rT}, \bar{\mathcal{X}}_{rT}) dw_{rt} \right] \\
&\quad dF^{sp}(\bar{\mathcal{Y}}_{rT}, \bar{\mathcal{X}}_{rT}, \bar{\mathbf{w}}_{r(t-1)}) \\
&= \frac{1}{T} \sum_{t=1}^T \int_{\bar{\mathcal{Y}}_{rT}, \bar{\mathcal{X}}_{rT}} \int_{w_{r1}} \int_{w_{r2}} \cdots \int_{w_{r(t-1)}} \left[\int_{w_{rt}} \frac{f_h(w_{rt})}{e_{rt}(w_{rt})} N_r(Y_{rt}(\bar{\mathbf{w}}_{rt})) f_{W_{rt}}(w_{rt} | \bar{H}_{r(t-1)}) dw_{rt} \right] \\
&\quad dF^{sp}(\bar{\mathcal{Y}}_{rT}, \bar{\mathcal{X}}_{rT}, \bar{\mathbf{w}}_{r(t-1)}) \quad (\text{From Assumption A.4}) \\
&= \frac{1}{T} \sum_{t=1}^T \int_{\bar{\mathcal{Y}}_{rT}, \bar{\mathcal{X}}_{rT}} \int_{w_{r1}} \int_{w_{r2}} \cdots \int_{w_{r(t-1)}} \left[\int_{w_{rt}} f_h(w_{rt}) N_r(Y_{rt}(\bar{\mathbf{w}}_{rt})) dw_{rt} \right] dF^{sp}(\bar{\mathcal{Y}}_{rT}, \bar{\mathcal{X}}_{rT}, \bar{\mathbf{w}}_{r(t-1)}) \\
&\quad (\text{From the definition of the region-specific propensity score}) \\
&= \frac{1}{T} \sum_{t=1}^T E_{F^{sp}} [N_{rt}(F_h)] \\
&= E_{F^{sp}} [N_r(F_h)]
\end{aligned}$$

Consistency We use an alteration of Lemma A in Section 7.2.1 of [Serfling \(1980\)](#). Since $\psi_r(\mathcal{D}_r; \mu)$ is monotone in μ with $\partial \psi_r(\mathcal{D}_r; \mu) / \partial \mu = -1 < 0$, we have that $\Psi_R(\mathcal{D}; \mu)$ and $\Psi^{sp}(\mu)$ are also monotone which implies uniqueness of their roots, $\hat{\mu}$ and μ_0 . From the strong law of large numbers we have that $\Psi_R(\mathcal{D}; \mu) \xrightarrow{a.s.} \Psi^{sp}(\mu)$, and

$$|\Psi^{sp}(\hat{\mu}) - \Psi^{sp}(\mu_0)| = |\Psi^{sp}(\hat{\mu}) - \Psi_R(\hat{\mu})| \leq \sup_{\mu} |\Psi^{sp}(\mu) - \Psi_R(\mu)| \rightarrow 0,$$

which, by the uniqueness of the roots for Ψ^{sp} and Ψ_R implies that $\hat{N}(F_h) \xrightarrow{a.s.} N(F_h)$ and $\hat{N}(F_h)$ is consistent for $N(F_h)$.

Asymptotic normality For asymptotic normality we will use Theorem A in Section 7.2.2 of [Serfling \(1980\)](#). We have already shown that μ_0 is an isolated root of $\Psi^{sp}(\mu) = 0$ (since it is unique) and that $\psi_r(\mathcal{D}_r; \mu)$ is monotone in μ . We also have that $\Psi^{sp}(\mu)$ is differentiable in μ with $\frac{\partial}{\partial \mu} \Psi^{sp}(\mu) = -1 \neq 0$. Lastly we will show that $E_{F^{sp}} [\psi_r^2(\mathcal{D}_r; \mu)]$ is finite in a neighborhood of μ_0 . To do so, consider μ in an ϵ -neighborhood of μ_0 , $\mu \in (\mu_0 - \epsilon, \mu_0 + \epsilon)$. Then

$$\begin{aligned}
E_{F^{sp}} [\psi_r^2(\mathcal{D}_r; \mu)] &= E_{F^{sp}} \left\{ \left[\frac{1}{T} \sum_{t=1}^T \frac{f_h(W_{rt})}{e_{rt}(W_{rt})} N_r(Y_{rt}) - \mu \right]^2 \right\} \\
&= E_{F^{sp}} \left\{ \left| \frac{1}{T} \sum_{t=1}^T \frac{f_h(W_{rt})}{e_{rt}(W_{rt})} N_r(Y_{rt}) - \mu \right|^2 \right\} \\
&\leq E_{F^{sp}} \left\{ \left[\frac{1}{T} \sum_{t=1}^T \frac{f_h(W_{rt})}{e_{rt}(W_{rt})} N_r(Y_{rt}) + |\mu| \right]^2 \right\} && \text{(Triangle inequality)} \\
&= E_{F^{sp}} \left\{ \left[\frac{1}{T} \sum_{t=1}^T \frac{f_h(W_{rt})}{e_{rt}(W_{rt})} N_r(Y_{rt}) \right]^2 \right\} + 2|\mu| E_{F^{sp}} \left[\frac{1}{T} \sum_{t=1}^T \frac{f_h(W_{rt})}{e_{rt}(W_{rt})} N_r(Y_{rt}) \right] + \mu^2 \\
&\leq E_{F^{sp}} \left\{ \left[\frac{1}{T} \sum_{t=1}^T \delta_W \delta_Y \right]^2 \right\} + 2|\mu| E_{F^{sp}} \left[\frac{1}{T} \sum_{t=1}^T \delta_W \delta_Y \right] + |\mu|^2 \\
&= (\delta_W \delta_Y)^2 + 2|\mu| \delta_W \delta_Y + \mu^2
\end{aligned}$$

where we used that all terms in the summation are positive along with Assumptions [A.6](#) and [A.7](#). Since $\mu \in (\mu_0 - \epsilon, \mu_0 + \epsilon)$ it is bounded, so the expectation above exists.

Then, since all the conditions of the theorem are satisfied we have that

$$\sqrt{R} \left(\widehat{N}(F_h) - N(F_h) \right) \rightarrow N(0, \sigma^2),$$

where $\sigma^2 = E_{F^{sp}} [\psi_r^2(\mathcal{D}_r; \mu_0)]$. □

C The Hájek Estimator

The standardization of weights used in the Hájek estimator is known to be effective in the settings where the weights are extreme. Its sample boundedness property guarantees that the resulting estimate is always within the range of the observed outcome. In our case, the Hájek estimator replaces the division by $T - M + 1$ with that by $\sum_{t=M}^T w_t$ where w_t is the product of fractions in Equation [\(6\)](#). For example,

$$\widehat{N}_B(F_h^M)_{\text{Hájek}} = \frac{1}{\sum_{t=M}^T w_t} \sum_{t=M}^T \widehat{N}_{Bt}(F_h^M)$$

The new martingale theorem stated in Theorem A.3 can be used in future research to show that the Hájek estimator is consistent and asymptotically normal, and derive the functional form of its asymptotic variance. However, for now, we use a heuristic approach to estimating the variance bound of the Hájek estimator. Since the Hájek estimator simply rescales the corresponding IPW estimator by $(T - M + 1) / \sum_{t=M}^T w_t$, we scale the variance bound derived for the estimator by $[(T - M + 1) / (\sum_{t=M}^T w_t)]^2$.

D Sensitivity analysis

In this section we discuss sensitivity analysis for the IPW estimators. In the main text of the manuscript we discuss sensitivity analysis for the Hájek estimator, which is admittedly a much harder problem due to the standardization of weights performed in the Hájek correction.

In this section we discuss sensitivity analysis based on the IPW estimator. We quickly see that bounding the estimator for different amounts of propensity score misspecification Γ can be directly achieved by solving a linear program. We can similarly bound the causal effect estimator exactly. In contrast, in the main text, bounding the value of the Hájek estimator requires additional tools to *transform* the problem to a linear program. This transformation forbids us from bounding the effect estimator exactly and forces us to acquire possibly conservative bounds for the effect estimator (see Theorem 4).

D.1 For the IPW estimator

We focus again on bounding the estimators for intervention over a single time period, though extensions to multiple time periods are direct, and discussed in more detail for the Hájek estimator in Appendix D.2. The IPW estimators that use the correct propensity score can be written as:

$$\begin{aligned}\widehat{N}_\rho(F_h) &= \frac{1}{T} \sum_{t=1}^T \rho_t w_t(F_h) \widetilde{N}_B(Y_t), \quad \text{and} \\ \widehat{\tau}_\rho(F_{h_1}, F_{h_2}) &= \frac{1}{T} \sum_{t=1}^T \rho_t w_t(F_{h_2}) \widetilde{N}_B(Y_t) - \frac{1}{T} \sum_{t=1}^T \rho_t w_t(F_{h_1}) \widetilde{N}_B(Y_t) \\ &= \frac{1}{T} \sum_{t=1}^T \rho_t [w_t(F_{h_2}) - w_t(F_{h_1})] \widetilde{N}_B(Y_t)\end{aligned}$$

where

$$w_t(F_h) = \frac{f_h(W_t)}{e_t(W_t)} \quad \text{and} \quad \widetilde{N}_B(Y_t) = \int_B \sum_{s \in S_{Y_t}} K_b(\|\omega - s\|) d\omega.$$

Both of the IPW estimators $\widehat{N}_\rho(F_h)$ and $\widehat{\tau}_\rho(F_{h_1}, F_{h_2})$ are linear in ρ , so finding their maximum/minimum over $\rho_t \in [\Gamma^{-1}, \Gamma]^T$ for each t is a linear problem and can be easily solved.

D.2 For the Hájek estimator

The standardization of the weights in the Hájek estimator implies that maximizing/minimizing the value of the estimator is no longer linear in ρ_t . This is evident in the form of the Hájek estimator for the number of points and the effect in a region, defined respectively as

$$\frac{\sum_{t=1}^T \rho_t w_t(F_h) \tilde{N}_B(Y_t)}{\sum_{t=1}^T \rho_t w_t(F_h)} \quad \text{and} \quad \frac{\sum_{t=1}^T \rho_t w_t(F_{h_2}) \tilde{N}_B(Y_t)}{\sum_{t=1}^T \rho_t w_t(F_{h_2})} - \frac{\sum_{t=1}^T \rho_t w_t(F_{h_1}) \tilde{N}_B(Y_t)}{\sum_{t=1}^T \rho_t w_t(F_{h_1})}$$

where

$$w_t(F_h) = \frac{f_h(W_t)}{e_t(W_t)} \quad \text{and} \quad \tilde{N}_B(Y_t) = \int_B \sum_{s \in S_{Y_t}} K_b(\|\omega - s\|) d\omega.$$

Theorem 4 states that bounding the estimator for the expected number of points $\hat{N}_\rho(F_h)$ can be transformed to a linear problem. However, the standardization of weights in the Hájek estimator and the fact that our estimator is the difference of two linear fractionals forbids us to see the problem of bounding the effect estimator $\hat{\tau}_\rho(F_{h_1}, F_{h_2})$ the same way.

Proof of Theorem 4. We view the problem of bounding $\hat{N}_\rho(F_{h_2})$ as a maximization/minimization problem of a linear fractional with positive denominator. These problems have been previously studied, and it has been shown that they can be transformed to a linear programming problem using the Charnes-Cooper transformation (Charnes and Cooper, 1962). The theorem states this transformation in the context of our estimator.

For the problem of bounding the effect estimator, the objective can be written as,

$$\begin{aligned} \hat{\tau}_\rho(F_{h_1}, F_{h_2}) &= \frac{\sum_{t=1}^T \rho_t w_t(F_{h_2}) \tilde{N}_B(Y_t)}{\sum_{t=1}^T \rho_t w_t(F_{h_2})} - \frac{\sum_{t=1}^T \rho_t w_t(F_{h_1}) \tilde{N}_B(Y_t)}{\sum_{t=1}^T \rho_t w_t(F_{h_1})} \\ &= \hat{N}_\rho(F_{h_2}) - \hat{N}_\rho(F_{h_1}). \end{aligned}$$

Thus, maximizing $\hat{\tau}_\rho(F_{h_1}, F_{h_2})$ over $\rho \in [\Gamma^{-1}, \Gamma]^T$ is equivalent to maximizing $\hat{N}_\rho(F_{h_2}) - \hat{N}_\rho(F_{h_1})$ over the same region for ρ . Since the space $(\rho_1, \rho_2) \in [\Gamma^{-1}, \Gamma]^{2T}$ includes (ρ_1, ρ_2) where $\rho_1 = \rho_2$ as a subspace, we have that

$$\begin{aligned} \max_{\rho \in [\Gamma^{-1}, \Gamma]^T} \left\{ \hat{N}_\rho(F_{h_2}) - \hat{N}_\rho(F_{h_1}) \right\} &\leq \max_{(\rho_1, \rho_2) \in [\Gamma^{-1}, \Gamma]^{2T}} \left\{ \hat{N}_{\rho_2}(F_{h_2}) - \hat{N}_{\rho_1}(F_{h_1}) \right\} \\ &= \max_{\rho_2 \in [\Gamma^{-1}, \Gamma]^T} \left\{ \hat{N}_{\rho_2}(F_{h_2}) \right\} - \min_{\rho_1 \in [\Gamma^{-1}, \Gamma]^T} \left\{ \hat{N}_{\rho_1}(F_{h_1}) \right\}, \end{aligned}$$

where the last equality holds since $\hat{N}_{\rho_j}(F_{h_j}) \geq 0$. Similarly, we can derive the bound for the minimum of $\hat{\tau}_\rho(F_{h_1}, F_{h_2})$. \square

Next, we derive similar conservative bounds for the estimators corresponding to the interventions over multiple time periods. Recall that the Hájek estimator for the number of events in region B under a stochastic intervention is given by,

$$\frac{\sum_{t=M}^T w_t(F_h^M) \tilde{N}_B(Y_t)}{\sum_{t=M}^T w_t(F_h^M)}, \quad \text{where} \quad w_t(F_h^M) = \prod_{j=t-M+1}^t \frac{f_h(W_j)}{e_j(W_j)}.$$

So, our sensitivity analysis would search to find the bounds of

$$\hat{N}_\rho(F_h^M) = \frac{\sum_{t=M}^T \left(\prod_{j=t-M+1}^t \rho_j \right) w_t(F_h^M) \tilde{N}_B(Y_t)}{\sum_{t=M}^T \left(\prod_{j=t-M+1}^t \rho_j \right) w_t(F_h^M)}$$

over $\rho \in [\Gamma^{-1}, \Gamma]$. Since each ρ_t can take a value in $[\Gamma^{-1}, \Gamma]$, the sensitivity analysis weights in the Hájek estimator for multiple time periods, $\prod_{j=t-M+1}^t \rho_j$, take a value in $[\Gamma^{-M}, \Gamma^M]$. Therefore the set $\{\alpha \in [\Gamma^{-M}, \Gamma^M]\}$ includes all vectors of length T whose t^{th} entry can be written as $\prod_{j=t-M+1}^t \rho_j$ for some vector ρ . Using the argument similar to that of Theorem 4, we have

$$\begin{aligned} \min_{\rho \in [\Gamma^{-1}, \Gamma]} \left\{ \hat{N}_\rho(F_h^M) \right\} &\geq \min_{\alpha \in [\Gamma^{-M}, \Gamma^M]} \left\{ \frac{\sum_{t=M}^T \alpha_t w_t(F_h^M) \tilde{N}_B(Y_t)}{\sum_{t=M}^T \alpha_t w_t(F_h^M)} \right\}, \\ \max_{\rho \in [\Gamma^{-1}, \Gamma]} \left\{ \hat{N}_\rho(F_h^M) \right\} &\leq \max_{\alpha \in [\Gamma^{-M}, \Gamma^M]} \left\{ \frac{\sum_{t=M}^T \alpha_t w_t(F_h^M) \tilde{N}_B(Y_t)}{\sum_{t=M}^T \alpha_t w_t(F_h^M)} \right\}. \end{aligned}$$

The quantities on the right can be computed by turning the linear fractional problem to a linear problem via the Charnes-Cooper transformation. Then, we can use these quantities as the conservative bounds for the minimum and maximum of our target quantities. Based on these bounds, we can again use Theorem 4 to acquire conservative bounds of the effect of changing the intervention for interventions over multiple time periods.

E Additional Simulation Results on the Iraq-based scenario

E.1 Asymptotic Variance and Bound, and Estimated Variance Bound

Figure 5 shows the average (over 200 simulated data sets) of the true asymptotic standard deviation and true bound as well as the estimated standard deviation bound of the IPW estimator for the average potential outcome using the true propensity score, for interventions taking place over $M \in \{1, 3\}$ time periods. Figure A.2 is a similar plot for the interventions taking place over $M = 1, 3,$ and 7 (rows) time periods, and observed time series of length $T = 200, 400, 500$ (columns). These plots show the median and interquartile range of the asymptotic standard deviation, true bound, and estimated bound over 200

simulated data sets.

We begin by focusing on low uncertainty scenarios, corresponding to the interventions taking place over $M = 1$ or 3 time periods with the distribution resembling the actual data generating mechanism. We think that the intervention distribution resembles the data generating mechanism in scenarios where the intervention intensity is close to 5, which is the average number of treatment-active locations for the data generating process. In these scenarios, the asymptotic variance bound is distinctly higher than the true asymptotic variance, indicating that the inference based on the true asymptotic bound would be conservative. We find that in these low uncertainty scenarios, the estimated bound is close to the true bound. For that reason, we would expect the confidence intervals for the IPW estimator based on the estimated bound to have a higher coverage probability than its nominal coverage (see Appendix E.2 for the coverage results).

In contrast, under high uncertainty scenarios such as the interventions over longer time periods, e.g., $M = 7$, the asymptotic standard deviation and theoretical bound are essentially indistinguishable. However, under these scenarios, the estimate of the theoretical bound tends to be biased downwards, suggesting that the confidence intervals for the IPW estimator based on the estimated bound would be anti-conservative. Furthermore, we expect it to take a longer time series in order for the estimated bound to converge to its theoretical value when the intervention takes place over a longer time period.

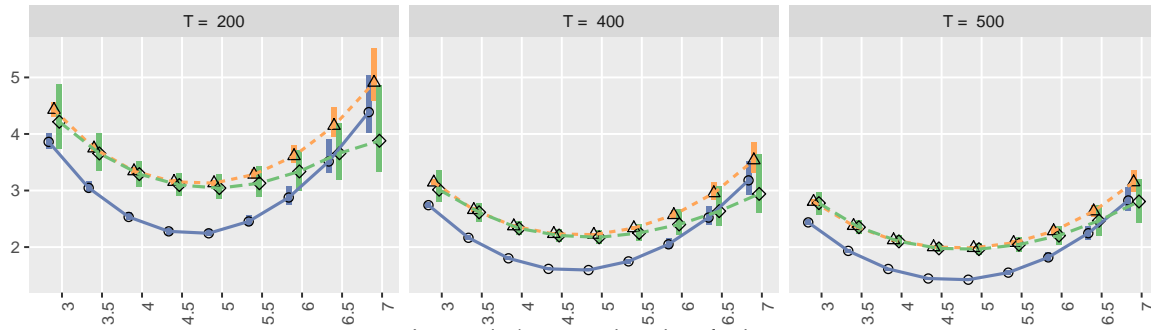
E.2 Coverage of the Confidence Intervals for the IPW and Hájek Estimators

IPW estimator. The results in Figure 5 indicate that the coverage of confidence intervals based on the asymptotic variance bound should be similar to those based on the true variance under high uncertainty scenarios, while they should be slightly higher under low uncertainty scenarios. Furthermore, confidence intervals based on the estimated variance bound should yield coverage probability close to (lower than) the coverage achieved using the theoretical bound under low (high) uncertainty scenarios.

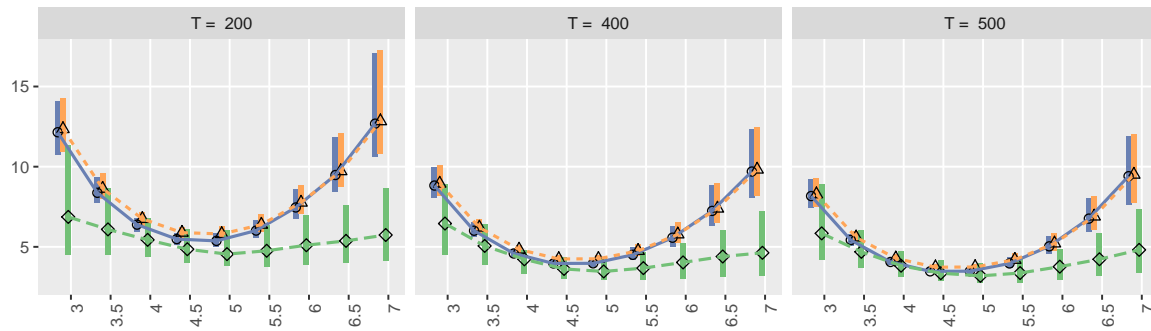
These expectations are indeed reflected in the coverage results shown in Figure A.3. Except when $M = 30$, the confidence interval for the IPW estimator based on either the *true* asymptotic variance or the *true* variance bound has a coverage of about 80% or higher. However, when $M = 30$, the confidence intervals based on the true asymptotic variance have a coverage below 60% or less, indicating that for interventions taking place over longer time periods, more data are needed to make use of the asymptotic approximation. However, these results are based on the true variance and variance bound, and instead inference would be based on the estimated variance bound. The under-estimation of the variance bound in high uncertainty scenarios found in Figure A.2 leads to the under-coverage of the confidence intervals based on the IPW estimator when using the estimated variance bound, especially when the interventions take place over long time periods.

Standard deviation and bound of the average potential outcome estimator for $B = \Omega$

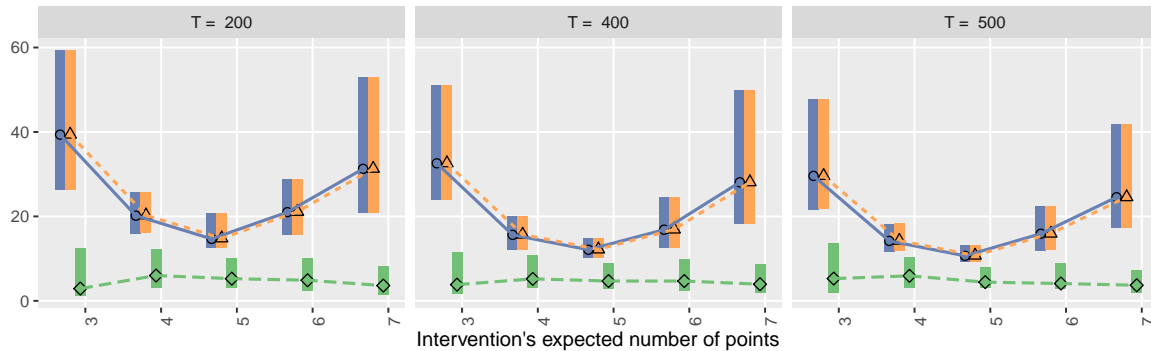
Intervention over 1 time point



Intervention over 3 time points



Intervention over 7 time points



○ True SD △ True SD Bound ◇ Estimated SD Bound

Figure A.2: Asymptotic Standard Deviation and Bound, and Estimated Bound. This figure shows the true asymptotic standard deviation (blue circles), the true asymptotic bound (orange triangles), and the estimated bound (green rhombuses) of the IPW estimator for the average potential outcome using the true propensity score, under interventions that take place over $M = 1, 3$ and 7 time periods (rows), and for increasing length of the time series (columns). The horizontal axis shows the intensity of the intervention at each time period. The points show the median value, and the rectangles show the interquartile range over 200 simulated data sets.

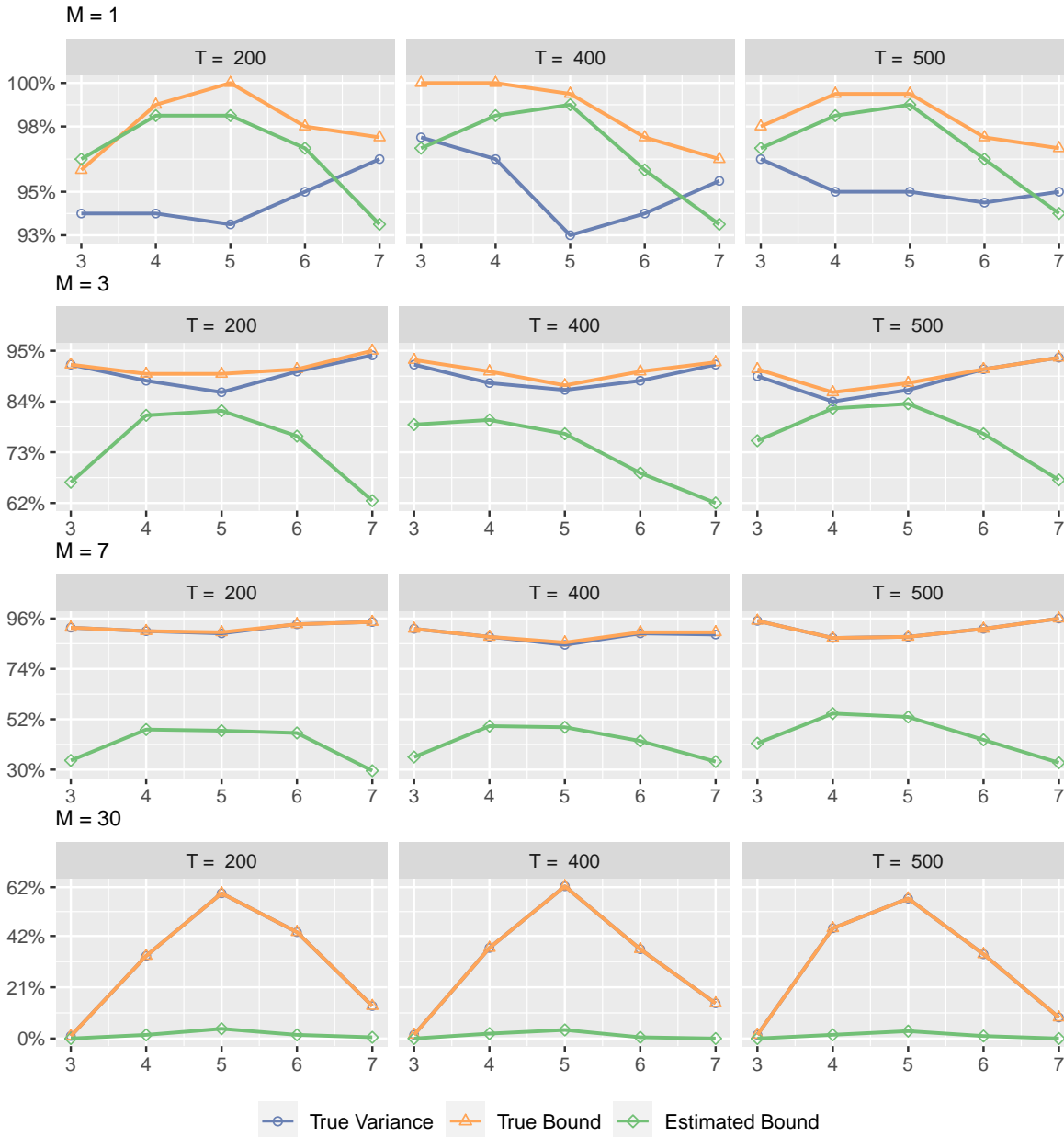


Figure A.3: Coverage of the IPW Estimator 95% Confidence Intervals. This figure shows the coverage of 95% confidence intervals for the average potential outcome over $B = \Omega$ based on the IPW estimator using the true variance (blue lines open circles), the true bound (orange lines with triangles), and the estimated bound (green lines with rhombuses), for interventions taking place over $M \in \{1, 3, 7, 30\}$ time periods (rows) and increasing length of the observed time series (columns).

Hájek estimator. Motivated by the good performance of the Hájek estimator shown in Figure 4, we also investigate the coverage probability of the 95% confidence interval as described in Appendix C. The rows of Figure A.4 show the coverage results for increasingly small regions, whereas the columns show the results for increasingly long observed time series ($T = 200, 400, 500$). Different colors correspond to the coverage results under interventions taking place over $M = 1$ (black), 3 (green), 7 (red), and 30

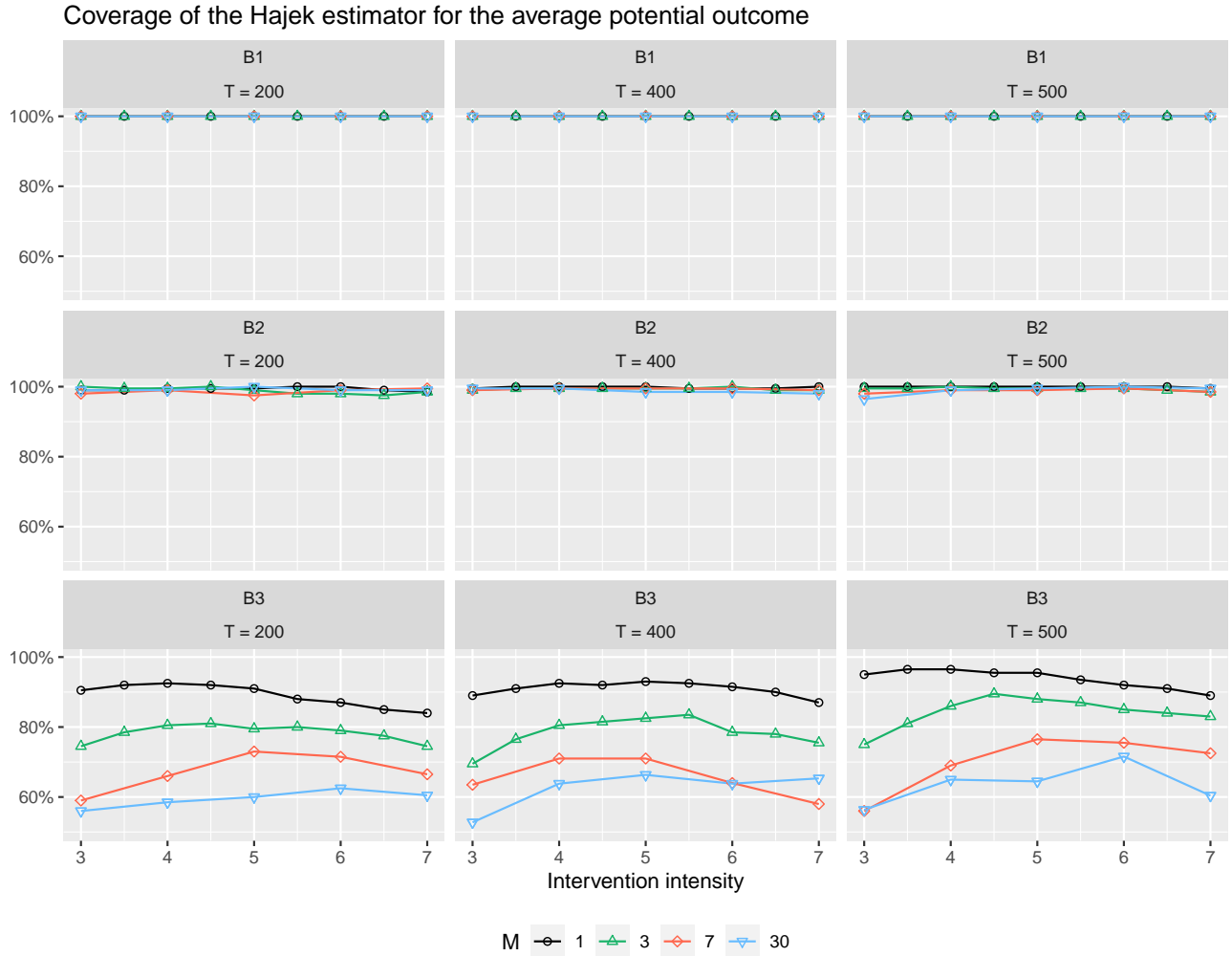


Figure A.4: Coverage of the Hájek Estimator’s 95% Confidence Intervals for the Average Potential Outcomes under Various Interventions. We vary the intervention intensity h (horizontal axis), and the length of intervention $M = 1, 3, 7, 30$ (different lines). Each row represents the coverage for different regions of interest, i.e., $B_1 = [0, 1]^2$, $B_2 = [0, 0.5]^2$ and $B_3 = [0.75, 1]^2$, whereas each column represents the length of time series, i.e., $T = 200, 400$ and 500 .

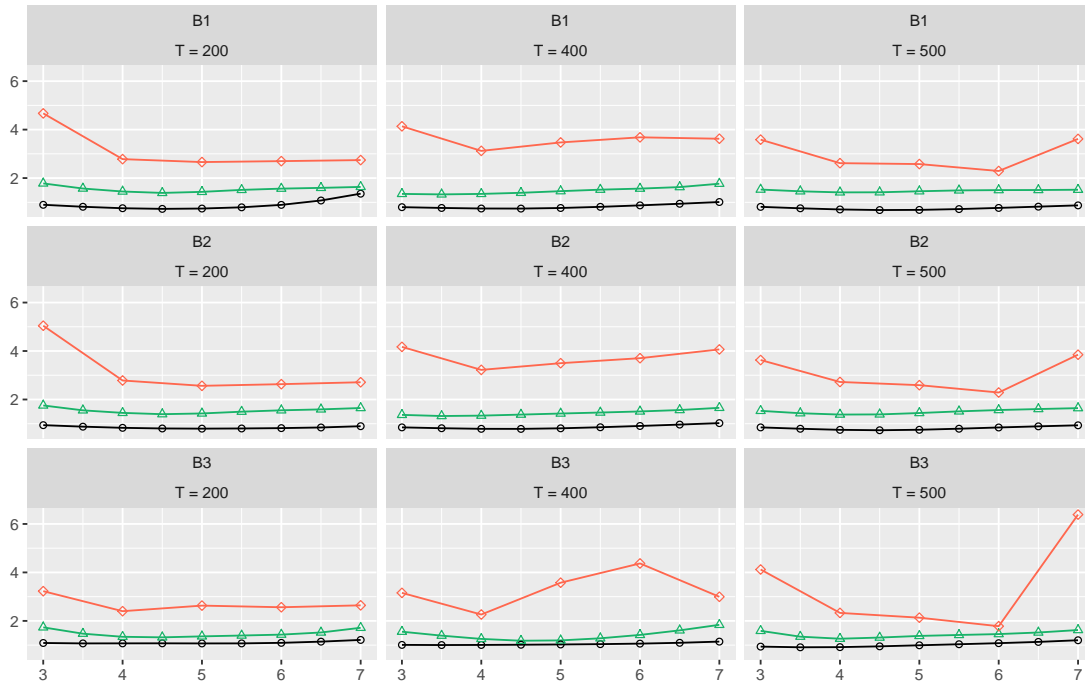
(blue) time periods. We find that the coverage is above 90% for all combinations of T and M for the two largest regions, even when an intervention takes place over 30 time periods. We find that the coverage is lower for the smallest region.

E.3 Uncertainty Estimates

Here, we show that the estimated standard deviation for the Hájek estimator outperforms that for the IPW estimator under many simulation scenarios.

We compute the standard deviation of the estimated average potential outcome across simulated data sets and compare it with the mean of the standard deviations, each of which is used to create the confidence intervals. The similarity of these two quantities implies the accuracy of our uncertainty

IPW estimator using the true propensity score



Hajek estimator

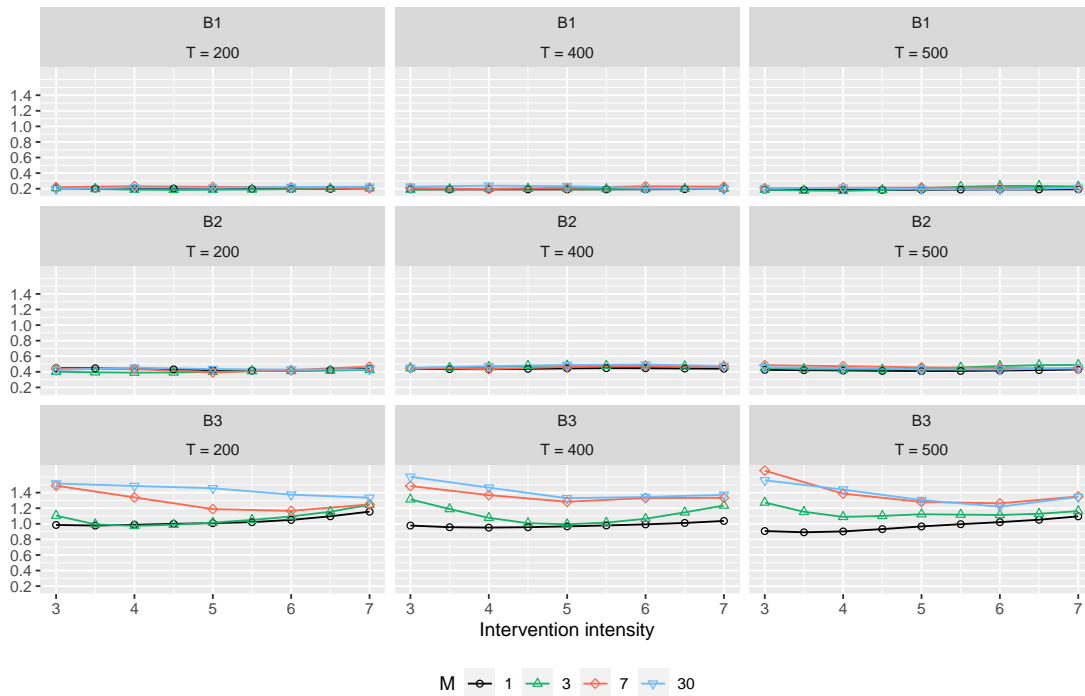


Figure A.5: Comparison of the Estimated and True Uncertainty for the Inverse Probability of Treatment and Hájek Estimators. Each plot presents the ratios between the standard deviation of each estimator and the mean estimated standard deviation across simulated data sets. A value smaller (greater) than 1 implies overestimation (underestimation) of uncertainty. The top (bottom) panel presents the results for the IPW (Hájek) estimator with the varying intensity under the intervention (horizontal axis) and for the whole country B_1 (first and fourth row) and two sub-regions, B_2 (second and fifth row) and B_3 (third and sixth row). We also vary the length of intervention, $M = 1, 3, 7$ and 30 time periods (black, green, red, and blue lines, respectively). The columns correspond to different lengths of the time series $T = 200, 400$ and 500 .

estimates. Figure A.5 presents the results as the ratio of these two quantities. A value below (above) 1 indicates that the true variability in our point estimates is smaller (greater) than our uncertainty estimate.

While the ratios are always below 1 for the Hájek estimator for the two largest regions B_1 and B_2 , they are almost always above 1 for the IPW estimator (top panel). This is consistent with the above results, showing that we tend to overestimate (underestimate) the uncertainty for the Hájek (IPW) estimator. We find that the confidence interval for the Hájek estimator tends to be most conservative when M is small and the region of interest is large. For the IPW estimator, the degree of uncertainty underestimation decreases as the length of time series T increases but increases as the length of intervention M increases. In fact, when $M = 30$, some of the ratios are as large as 20 (hence they are not included in the figure). The results suggest that in practice the Hájek estimator should be preferred over the IPW estimator especially for stochastic interventions over a long time period.

E.4 Covariate Balance

We evaluate the balance of covariates based on the estimated propensity score by comparing their p-values in the propensity score model, and in a model with functional form as in the propensity score model but weighted by the inverse of the estimated propensity score. The left plot of Figure A.6 shows the p-value for the previous outcome-active locations, which are one of the time-varying confounders, across 200 simulated data sets. Evidently, the p-values in the unweighted model are close to 0, indicating that previous outcome-active locations form an important predictor of the treatment assignment. However, in the weighted model, the p-values of the same confounder are more evenly distributed across the $(0, 1)$ range, indicating that this confounder is better balanced in the weighted time series.

F Additional simulations on a square geometry

F.1 The Simulation Design

We also consider a time series of point patterns of length $T \in \{200, 400, 500\}$ on the unit square, $\Omega = [0, 1] \times [0, 1]$. For each time series length T , 200 data sets are generated with the following design.

Time-varying and time-invariant confounders. Our simulation study includes two time-invariant and two time-varying confounders. For the first time-invariant confounder, we construct a hypothetical road network on Ω using lines and arcs, which is highlighted by bright white lines in Figure A.7a. Then, we define $X^1(\omega) = 1.2 \exp\{-2D_1(\omega)\}$ where $D_1(\omega)$ is the distance from ω to the closest line. The second time-invariant covariate is constructed similarly, as $X^2(\omega) = \exp\{-3D_2(\omega)\}$ where $D_2(\omega)$ is the distance to the closest arc. In addition, the time-varying confounders, $X_t^3(\omega)$ and $X_t^4(\omega)$, are defined based on the exponential decay of distance to the closest point; these points are generated according to a

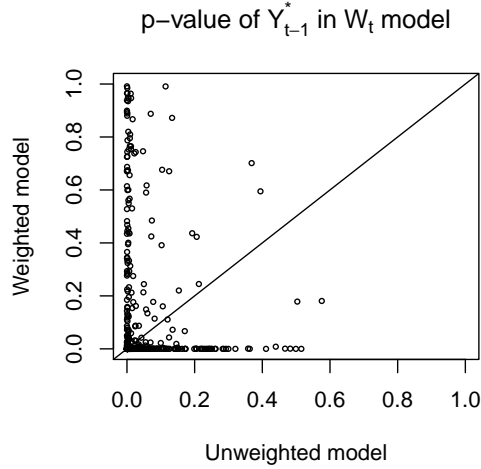


Figure A.6: Balance of the Previous Outcome-Active Locations in Treatment Model. Each point shows the relative magnitude of the p-value for the previous outcome-active locations in the unweighted propensity score model (horizontal axis) over that of the model weighted by the inverse of the estimated propensity score (vertical axis).

non-homogeneous Poisson point processes with the following intensity function

$$\lambda_t^{X^j}(\omega) = \exp \{ \rho_0^j + \rho_1^j X^1(\omega) \}, \quad j = 3, 4,$$

where $\rho_1^3 = 1$, and $\rho_1^4 = 1.5$. Figure A.7b shows one realization of $X_t^3(\omega)$.

Spatio-temporal point processes for treatment and outcome variables. We again generate treatment and outcome point patterns from non-homogeneous Poisson processes that depends on all confounders, and the previous treatment and outcome realizations. The functional specification of the Poisson process intensities is the same as in Section 6. The model gives rise to an average of 5 observed treatment-active locations and 21 observed outcome-active locations within each time period.

Stochastic interventions. We consider interventions of the form F_h^M based on a homogeneous Poisson process with intensity h that is constant over Ω and ranges from 3 to 7. We consider various lengths of each intervention by setting $M \in \{1, 3, 7, 30\}$. The second intervention we consider is defined over the three time periods, i.e., $F_h = F_{h_3} \times F_{h_2} \times F_{h_1}$ with $M = 3$. The intervention for the first time period F_{h_3} is a homogeneous Poisson process with intensity h_3 ranging from 3 to 7, whereas $F_{h_2} = F_{h_1}$ is a homogeneous Poisson process with intensity equal to 5 everywhere over Ω . For each stochastic intervention, we consider the region of interest, denoted by set B , of three different sizes: $B = \Omega = [0, 1] \times [0, 1]$, $B = [0, 0.5] \times [0, 0.5]$, and $B = [0.75, 1] \times [0.75, 1]$.

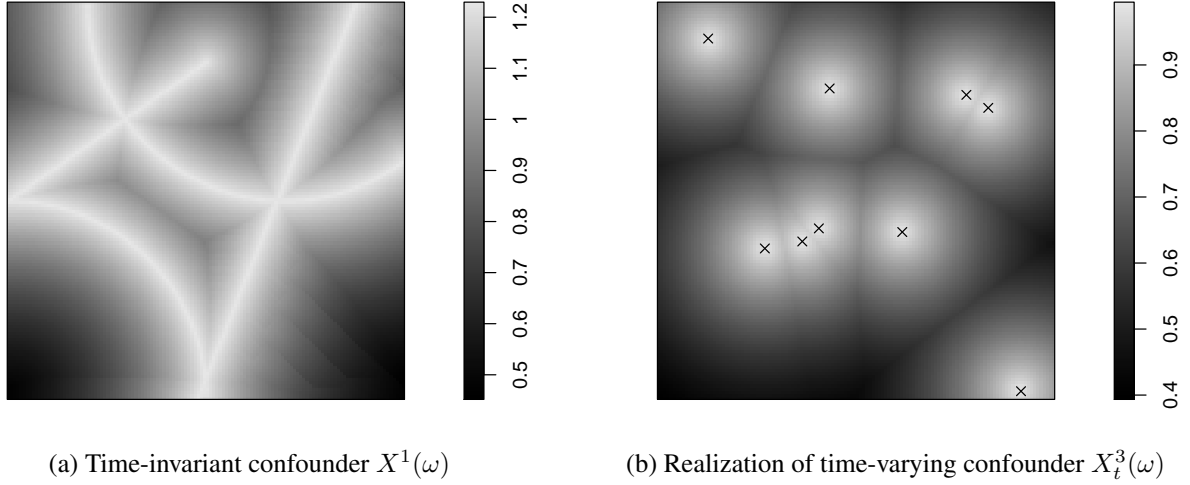


Figure A.7: Simulated Confounders. Panel (a) shows one of the two time-invariant confounders representing the exponential decay of distance to the road network. Panel (b) shows one realization for one of the time-varying confounders. Points \times are generated from a non-homogeneous Poisson process depending on the road network in (a). Then, the time-varying confounder is defined as the exponential decay of distance to the points \times .

Estimand and estimation. Approximating the true values of the estimands and estimation is performed as described in Section 6. In these simulations, for $T = 500$ (the longest time series in our simulation scenario) the spatial smoothing bandwidth is approximately equal to 0.16, smaller than the size of the smallest B (which is equal to $[0.75, 1]^2$).

Variance and its upper bound. We base calculation of the theoretical variance and the variance bound on Theorems 1 and A.1, and use Monte Carlo approximations to compute these, as in Section 6. We also use Lemma 1 We use Lemma 1 to estimate the variance bound.

Covariate balance. As in Section 6, we use weighted regression by the estimated propensity score to investigate covariate balance.

F.2 Simulation Results

Estimation. Figures A.8 and A.9 present the results. In Figure A.8, the top panel shows how the (true and estimated) average potential outcomes in the whole region ($B = \Omega$) change as the intensity varies under the single time period interventions. The bottom panel shows how the true and estimated average potential outcomes in the sub-region $[0.75, 1]^2$ change under the three time period interventions when the intensity at three time periods ago ranges from 3 to 7. For both simulation scenarios, we vary the length of the time series from 200 (left plots) to 500 (right plots).

As expected, the unadjusted estimates (green crosses) are far from the true average potential outcome (black solid circles) across all simulation scenarios. In contrast, and consistent with the results of

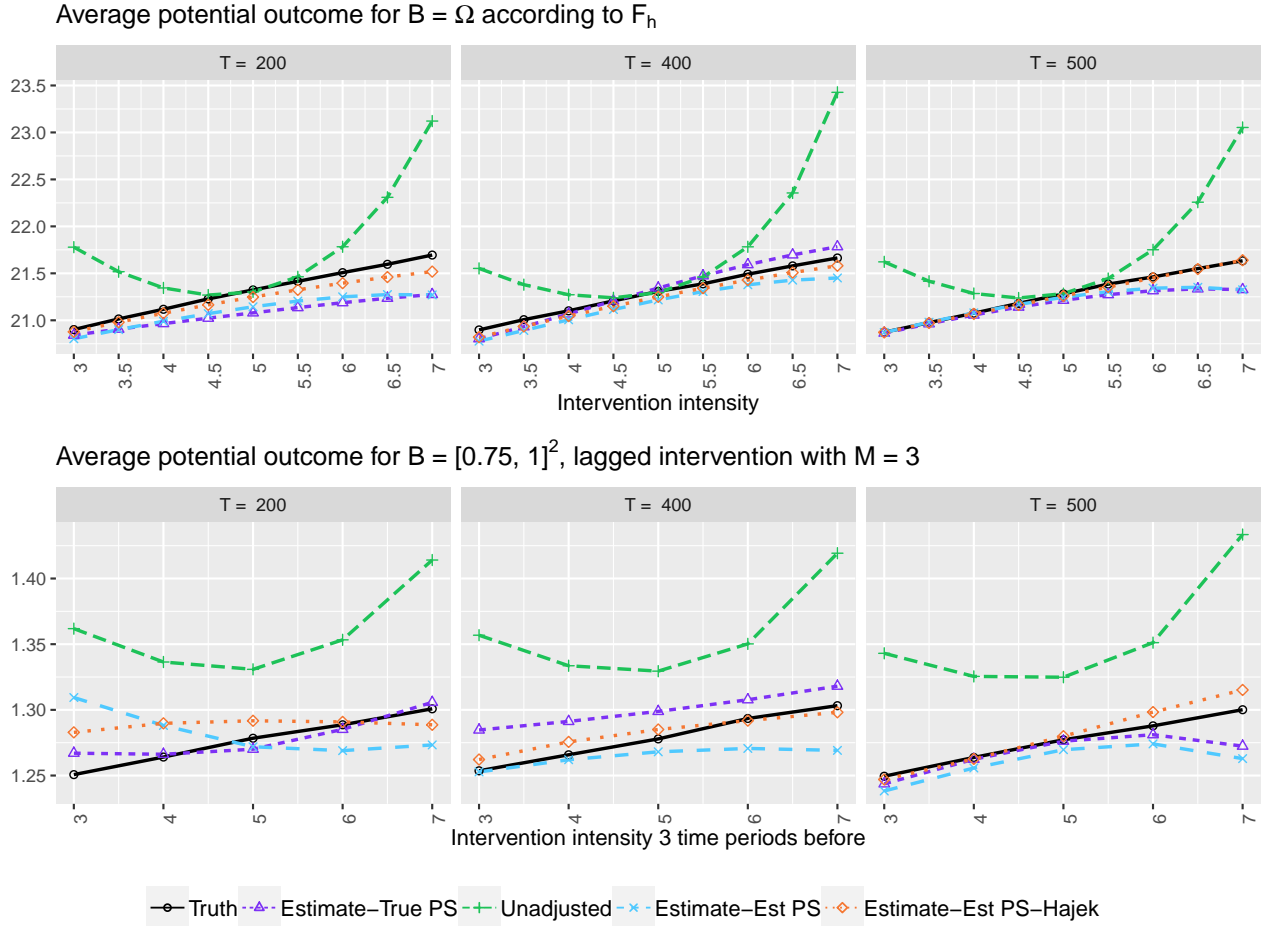


Figure A.8: Simulation Results for the True and Estimated Average Potential Outcomes. In the top panel, we present the true and estimated average potential outcomes in the entire region $B = \Omega$ under single-time interventions with the varying intensity (horizontal axis). In the bottom panel, we consider the average potential outcome in the sub-region $B = [0.75, 1]^2$ for the intervention F_h , with $M = 3$, the varying intensity of F_{h_3} (horizontal axis), and F_{h_1}, F_{h_2} intensity set to 5. The black lines with solid circles represent the truths, while the other dotted or dashed lines represent the estimates; the estimator based on the true propensity score (purple triangles), the unadjusted estimator (green crosses), the estimator based on the estimated propensity score (blue x's), the Hájek estimator based on the estimated propensity score (orange rhombuses).

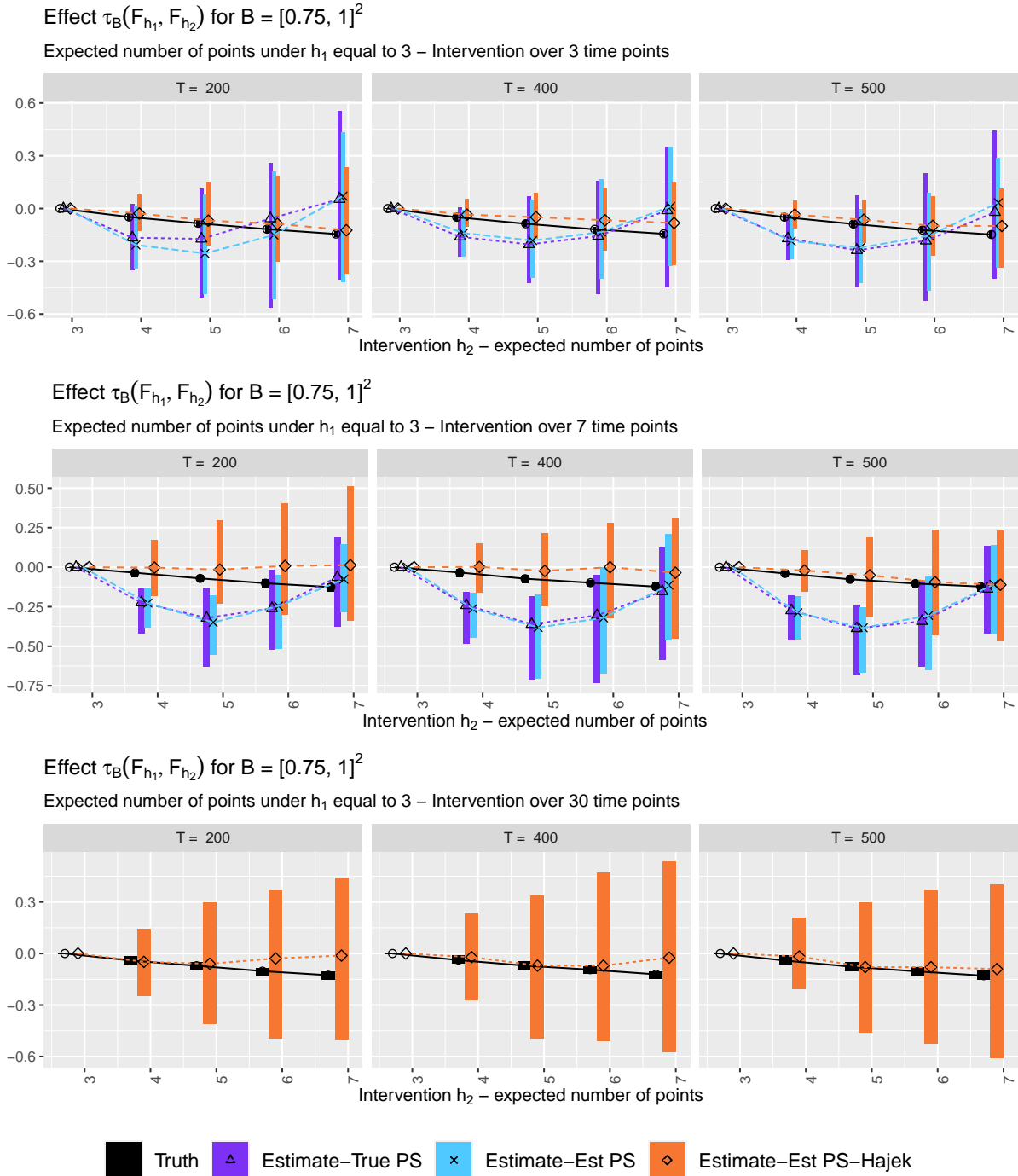


Figure A.9: Simulation Results for the Interventions of Increasing Time Lengths. Rows correspond to the interventions taking place over $M = 3, 7$, and 30 time periods. Columns correspond to the increasing length of the time series from 200 (left plots) to 500 (right plots). The vertical axis shows the change in the expected number of the outcome active locations over $[0.75, 1]^2$ for a change in the intervention intensity from 3 under h_1 to the value shown in the horizontal axis under h_2 , for M time periods. The points in the plot show the median estimate over 200 data sets, and the rectangles show the interquartile range of estimates. Only the Hájek estimates are shown for $M = 30$ as the extremely small weights arising from a large number of time periods make the estimates from the other estimators close to zero.

Theorems 1 and A.1, the accuracy of the proposed estimator (purple triangles based on the true propensity score, blue x's based on the estimated propensity score) improves as the number of time periods increases. We note that the convergence is slower when $M = 3$ than $M = 1$.

Figure A.9 shows the performance of the estimators for the interventions over many time periods. The plots show the estimated change in the number of outcome-active locations over the sub-region $B = [0.75, 1]$ for a change in the stochastic intervention from 3 per time period to the value on the horizontal axis. The rows correspond to the interventions over $M = 3, 7,$ and 30 time periods, respectively, whereas the columns represent the different lengths of time series, i.e., $T = 200, 400$ and 500 . The results are shown for the IPW estimators based on the true propensity score (purple lines with open triangles) and the estimated propensity score (blue lines with x's) as well as the Hájek estimator based on the estimated propensity score (orange lines with open rhombuses). Only the Hájek estimates are shown for $M = 30$ as the extremely small weights arising from a large number of time periods make the estimates from the other estimators essentially equal to zero. The lines and points in the plot show the median estimate and the rectangles show the interquartile range of estimates across 200 simulated data sets.

Again, as in the simulations of Section 6, we find that the Hájek estimator performs well across all simulation scenarios, whereas the IPW estimator tends to suffer from extreme weights.

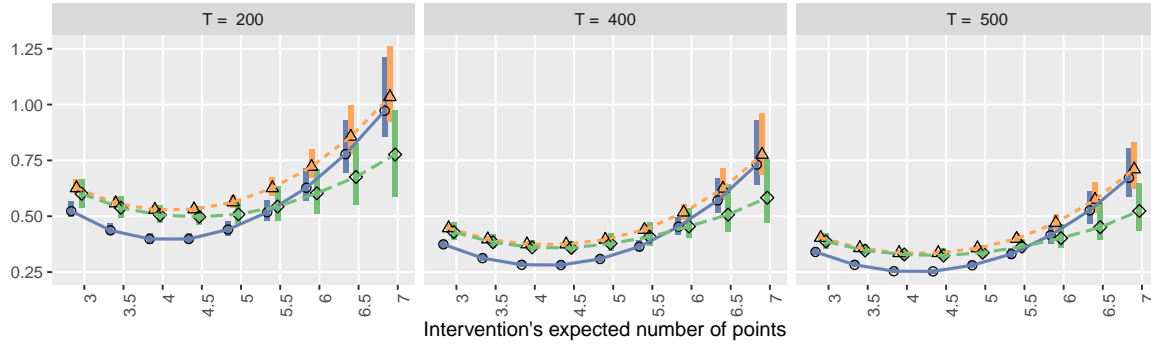
The variance and its bound. Next, we compare the true theoretical variance, v/T , with the variance bound v^*/T and its consistent estimator (see Lemma 1). We again focus on the proposed estimators with the true propensity score. Figure A.10 shows the results of an intervention F_h^M for $M = 1, 3$ and 7 , for region $B = [0, 0.5]^2$, and observed time series of length $T = 200, 400, 500$. These plots show the median and interquartile range of the asymptotic standard deviation, true bound, and estimated bound over 200 simulated data sets.

The conclusions are similar to the main manuscript. As expected, the true variance decreases as the total number of time periods increases. We start by focusing on low uncertainty scenarios, corresponding to the interventions taking place over $M = 1$ or 3 time periods with the distribution resembling the actual data generating mechanism. We think that the intervention distribution resembles the data generating mechanism in scenarios where the intervention intensity is close to 5, which is the average number of treatment-active locations for the data generating process. In these scenarios, the asymptotic variance bound is distinctly higher than the true asymptotic variance, indicating that the inference based on the true asymptotic bound would be conservative. We find that in these low uncertainty scenarios, the estimated bound is close to the true bound. For that reason, we would expect the confidence intervals for the IPW estimator based on the estimated bound to have a higher coverage probability than its nominal coverage.

In contrast, under high uncertainty scenarios such as the interventions over longer time periods, e.g., $M = 7$, the asymptotic standard deviation and theoretical bound are essentially indistinguishable. However, under these scenarios, the estimate of the theoretical bound tends to be biased downwards,

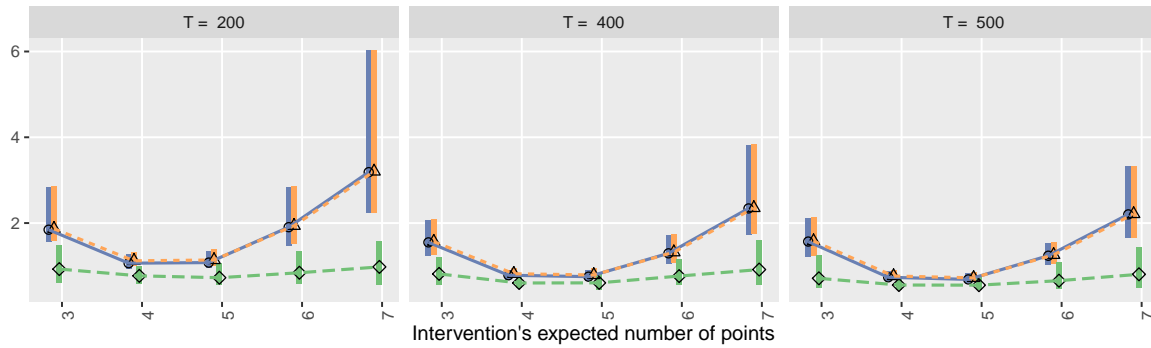
Standard deviation and bound of the average potential outcome estimator for $B = [0, 0.5]^2$

Intervention over 1 time point



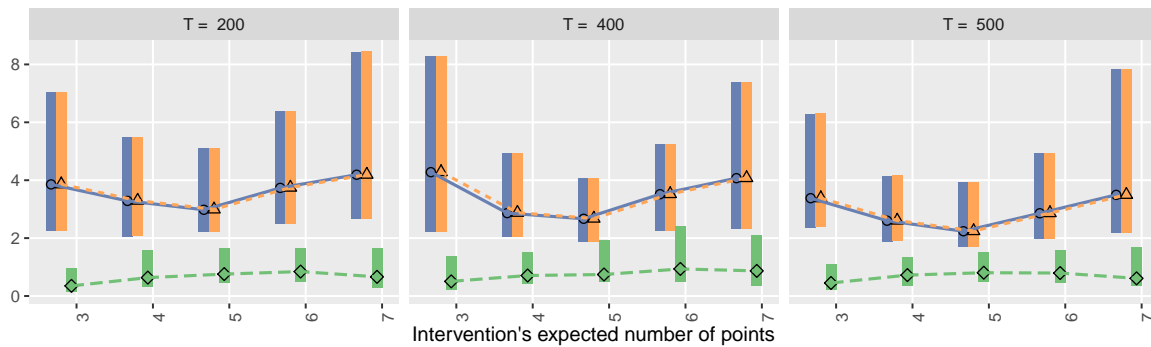
Standard deviation and bound of the average potential outcome estimator for $B = [0, 0.5]^2$

Intervention over 3 time points



Standard deviation and bound of the average potential outcome estimator for $B = [0, 0.5]^2$

Intervention over 7 time points



○ True SD △ True SD Bound ◇ Estimated SD Bound

Figure A.10: Asymptotic Standard Deviation and Bound, and Estimated Bound. This figure shows the true asymptotic standard deviation (blue circles), the true asymptotic bound (orange triangles), and the estimated bound (green rhombuses) of the IPW estimator for the average potential outcome using the true propensity score, under interventions that take place over $M = 1, 3$ and 7 time periods (rows), and for increasing length of the time series (columns). The horizontal axis shows the intensity of the intervention at each time period. The points show the median value, and the rectangles show the interquartile range over 200 simulated data sets.

suggesting that the confidence intervals for the IPW estimator based on the estimated bound would be anti-conservative. As the length of time series increases, the estimated variance bound more closely approximates its theoretical value (consistent with Lemma 1), but we expect it to take a longer time series in order for the estimated bound to converge to its theoretical value when the intervention takes place over a longer time period.

Coverage. These results on the asymptotic variance and variance bound lead to similar conclusions with respect to the coverage of 95% confidence intervals of the IPW estimator.

The coverage results are shown in Figure A.11. We find that, except when $M = 30$, the confidence interval for the IPW estimator based on either the true asymptotic variance or the true variance bound has a coverage of about 80% or higher. This implies that the asymptotic normality established in Theorem 1 provides an adequate approximation to the estimator’s sampling distribution for small or moderate values of M . However, for $M = 30$, the confidence interval for the IPW estimator is anti-conservative due to the fact that the weights, which equal the product of ratios across many time periods, become extremely small. In addition, the underestimation of the variance bound in high uncertainty scenarios found in Figure A.10 leads to the under-coverage of the confidence intervals based on the IPW estimator and using the estimated variance bound, especially when the interventions take place over long time periods.

We also investigate the coverage probability of the 95% confidence interval for the Hájek estimator. The rows of Figure A.12 show the coverage results for increasingly small regions, i.e., $B_1 = [0, 1]^2$, $B_2 = [0, 0.5]^2$, and $B_3 = [0.75, 1]^2$, whereas the columns show the results for increasingly long observed time series ($T = 200, 400, 500$). Different colors correspond to the coverage results under interventions taking place over $M = 1$ (black), 3 (green), 7 (red), and 30 (blue) time periods. We find that the coverage is above 85% for all cases, even when an intervention takes place over 30 time periods. As expected, the coverage is higher for smaller values of M , since these correspond to lower-uncertainty situations. We also find that the coverage is lower for smaller regions.

Comparison of Monte Carlo and estimated variance We find that the confidence interval for the Hájek estimator has a better coverage probability even for the interventions over long time periods. Here, we show that the estimated standard deviation for the Hájek estimator outperforms that for the IPW estimator under many simulation scenarios.

Figure A.13 shows the ratio of the standard deviation of the estimated average potential outcome across simulated data sets over the mean of the standard deviations. A value below (above) 1 indicates that the true variability in our point estimates is smaller (greater) than our uncertainty estimate. While the ratios are always below 1 for the Hájek estimator (bottom panel), they are almost always above 1 for the IPW estimator (top panel). This shows that we tend to overestimate (underestimate) the uncertainty for the Hájek (IPW) estimator. Further, we find that the confidence interval for the Hájek estimator tends to be most conservative when M is small and the region of interest is large. For the IPW estimator, the

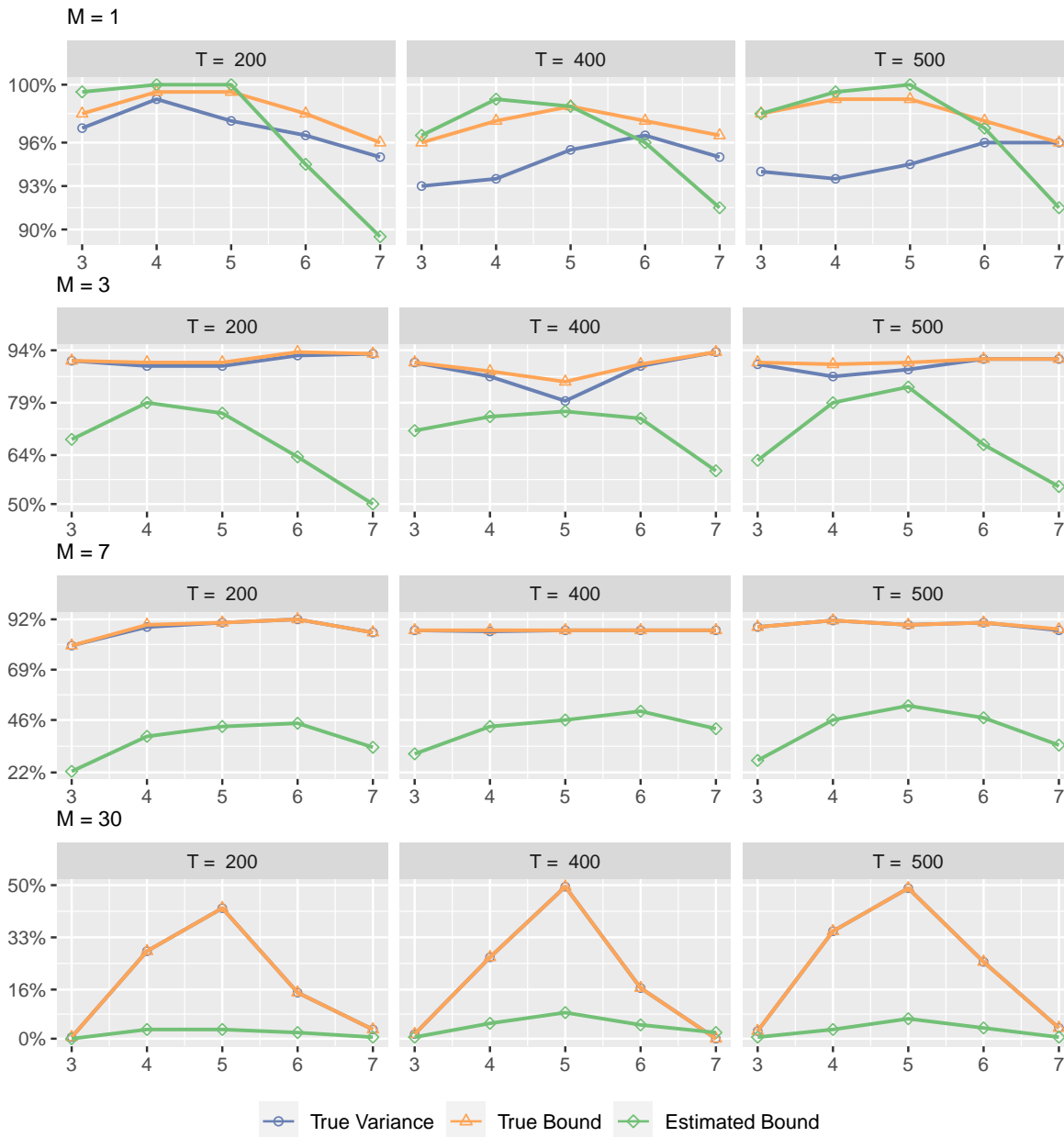


Figure A.11: Coverage of the IPW Estimator 95% Confidence Intervals. This figure shows the coverage of 95% confidence intervals for the average potential outcome over $B = \Omega$ based on the IPW estimator using the true variance (blue lines open circles), the true bound (orange lines with triangles), and the estimated bound (green lines with rhombuses), for interventions taking place over $M \in \{1, 3, 7, 30\}$ time periods (rows) and increasing length of the observed time series (columns).

degree of underestimation decreases as the length of time series T increases but increases as the length of intervention M increases. In fact, when $M = 30$, some of the ratios are as large as 20 (hence they are not included in the figure). The results suggest that in practice the Hájek estimator should be preferred over the IPW estimator especially for stochastic interventions over a long time period.

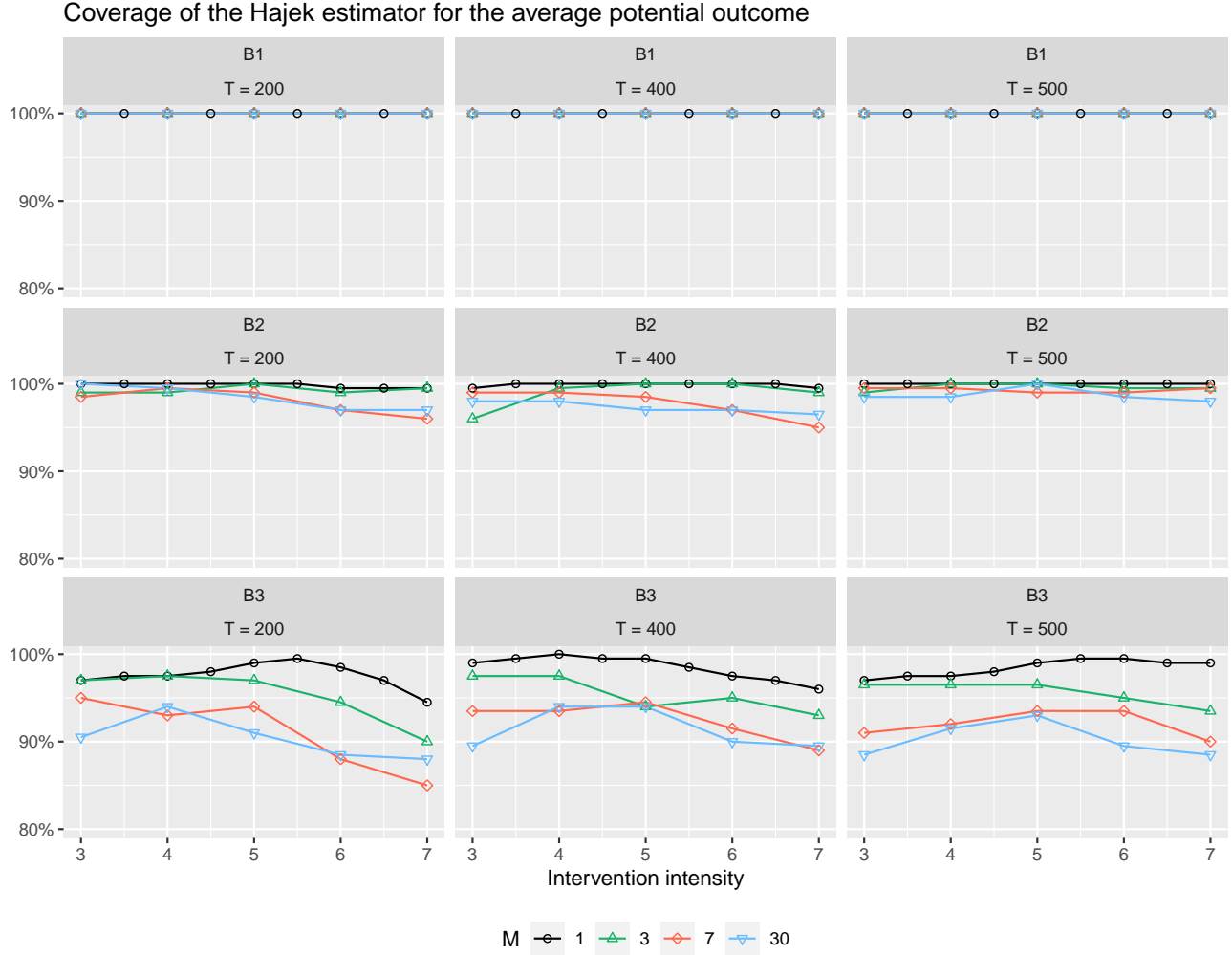


Figure A.12: Coverage of the Hájek Estimator’s 95% Confidence Intervals for the Average Potential Outcomes under Various Interventions. We vary the intervention intensity h (horizontal axis), and the length of intervention $M = 1, 3, 7, 30$ (different lines). Each row represents the coverage for different regions of interest, i.e., $B_1 = [0, 1]^2$, $B_2 = [0, 0.5]^2$ and $B_3 = [0.75, 1]^2$, whereas each column represents the length of time series, i.e., $T = 200, 400$ and 500 .

Balance. As in the simulations in the main manuscript, we find that the p-values of one of the confounders (Y_{t-1}^* in Equation(12)) are substantially greater in the weighted propensity score model than in the unweighted model, where the weights are given by the inverse of the estimated propensity score (shown in Figure A.14).

G Additional Empirical Results

G.1 Visualization

As discussed in Section 7.1, we consider a stochastic intervention whose focal point is the center of Baghdad. The degree of concentration is controlled by the precision parameter α whose greater value,

IPW estimator using the true propensity score

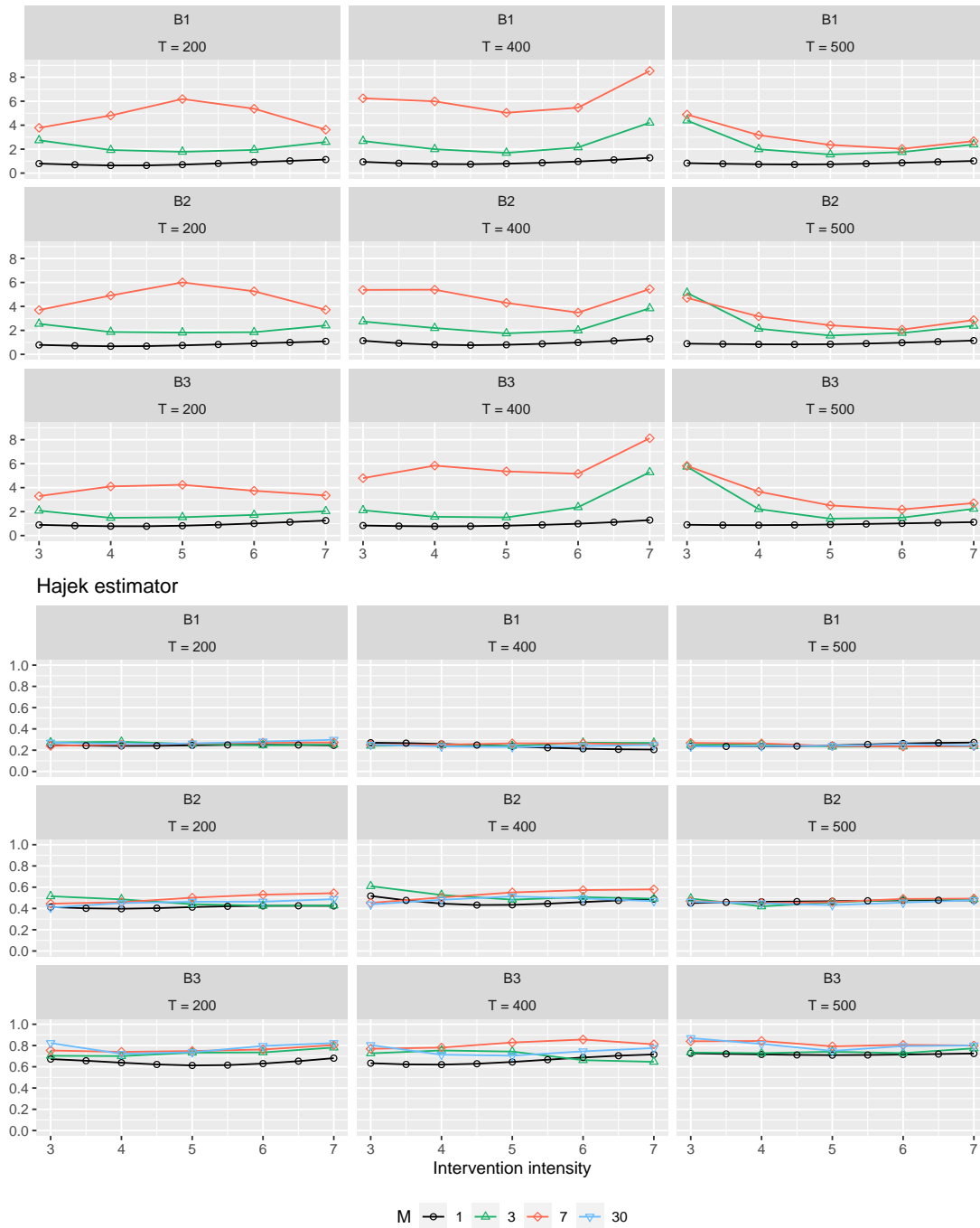


Figure A.13: Comparison of the Estimated and True Uncertainty for the Inverse Probability of Treatment and Hájek Estimators. Each plot presents the ratios between the standard deviation of each estimator and the mean estimated standard deviation across simulated data sets. A value smaller (greater) than 1 implies overestimation (underestimation) of uncertainty. The top (bottom) panel presents the results for the IPW (Hájek) estimator with the varying intensity under the intervention (horizontal axis) and for the whole region B_1 (first and fourth row) and two sub-regions, $B_2 = [0, 0.5]^2$ (second and fifth row) and $B_3 = [0.75, 1]^2$ (third and sixth row). We also vary the length of intervention, $M = 1, 3, 7$ and 30 time periods (black, green, red, and blue lines, respectively). The columns correspond to different lengths of the time series $T = 200, 400$ and 500.

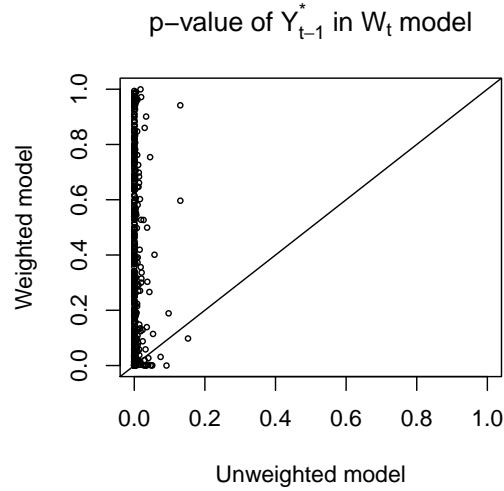


Figure A.14: Balance of the Previous Outcome-Active Locations in Treatment Model. In the left plot, each point shows the relative magnitude of the p-value for the previous outcome-active locations in the unweighted propensity score model (horizontal axis) over that of the model weighted by the inverse of the estimated propensity score (vertical axis). The right plot shows the distribution of the estimated coefficient of the previous outcome-active locations in the weighted propensity score model.

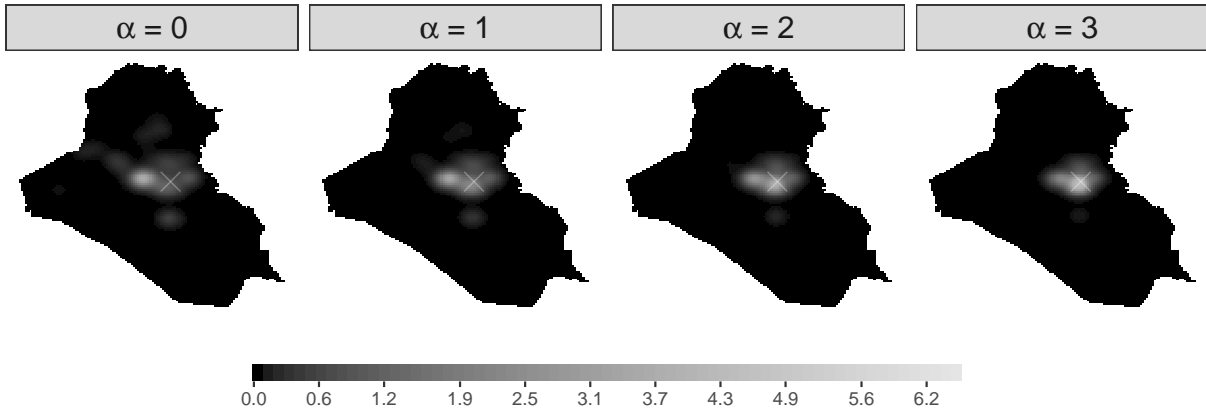


Figure A.15: Visualization of Intensity under Stochastic Interventions whose Focal Point is the Center of Baghdad. Across plots, we vary the degree to which the airstrikes are concentrated around the focal point using the precision parameter, while the expected number of airstrikes is held constant at 3 per day.

implying that more airstrikes are occurring near the focal point. We vary the value of α from 0 to 3, while keeping the expected number of airstrikes constant at 3 per day. Figure A.15 illustrates intensities for the different values of α . The first plot in the figure does not focus on Baghdad at all, representing the baseline spatial distribution ϕ_0 . As the value of α increases, the spatial distribution of airstrikes becomes concentrated more towards the center of Baghdad.

Type ($F_{h'}$, $F_{h''}$)	M	Outcome	Iraq	Baghdad	Outside Baghdad
Increasing the intensity (1, 3)	3	IED	-1.3 (-8.2, 5.7)	-0.2 (-2.2, 1.8)	-1.1 (-6.2, 4)
		SAF	-1.9 (-8.9, 5.2)	-1.2 (-3.8, 1.4)	-0.7 (-5.2, 3.9)
		Other Attack	-1.8 (-19.9, 16.4)	0.1 (-6.2, 6.5)	-1.9 (-13.9, 10)
	7	IED	5 (-2.6, 12.7)	1.2 (-0.8, 3.3)	3.8 (-1.9, 9.5)
		SAF	10 (1.7, 18.2)	3.1 (0.4, 5.8)	6.9 (1.1, 12.6)
		Other Attack	14 (-5.2, 33.3)	5.6 (-0.4, 11.6)	8.4 (-5.1, 21.9)
	30	IED	11 (-1.1, 23)	3.2 (-0.1, 6.5)	7.8 (-1.1, 16.6)
		SAF	14.8 (2.9, 26.7)	5.9 (1.2, 10.5)	8.9 (1.5, 16.3)
		Other Attack	33.1 (3.3, 62.9)	13.6 (2.7, 24.6)	19.5 (0.3, 38.6)
Changing the focal points (0, 3)	3	IED	2.6 (-7.2, 12.5)	0.6 (-2.4, 3.5)	2.1 (-4.9, 9)
		SAF	2.6 (-6.7, 12)	0.9 (-2.2, 4.1)	1.7 (-4.8, 8.1)
		Other Attack	7.5 (-16, 31)	3.8 (-4.3, 11.8)	3.7 (-12, 19.4)
	7	IED	2 (-6.9, 10.8)	1.1 (-1.7, 3.9)	0.9 (-5.2, 7)
		SAF	0.2 (-9.7, 10.1)	0.6 (-2.8, 3.9)	-0.4 (-7.1, 6.4)
		Other Attack	3.5 (-18.8, 25.8)	1.7 (-6, 9.4)	1.8 (-13, 16.7)
	30	IED	-1.2 (-15.9, 13.4)	-0.7 (-4.5, 3.1)	-0.5 (-11.6, 10.6)
		SAF	5.7 (-10, 21.4)	-1.3 (-6.4, 3.8)	7 (-4.1, 18.1)
		Other Attack	-3.5 (-37.8, 30.7)	-6.6 (-17.5, 4.2)	3.1 (-21.1, 27.3)
Lagged effects (1, 5)	3	IED	-2.3 (-10.1, 5.5)	-0.6 (-2.7, 1.5)	-1.7 (-7.5, 4.1)
		SAF	-1 (-9.9, 8)	-0.7 (-4.5, 3)	-0.2 (-5.5, 5)
		Other Attack	-3.9 (-23.6, 15.8)	-1.2 (-8.2, 5.9)	-2.8 (-15.5, 10)
	7	IED	6.8 (-0.7, 14.3)	2.2 (-0.2, 4.6)	4.6 (-0.6, 9.8)
		SAF	9.4 (1.6, 17.2)	3.6 (1, 6.2)	5.8 (0.4, 11.2)
		Other Attack	20.9 (2.3, 39.4)	8.2 (1.8, 14.6)	12.7 (0.4, 24.9)
	30	IED	1.5 (-3.8, 6.8)	0.3 (-1, 1.5)	1.2 (-2.8, 5.3)
		SAF	2.8 (-1.8, 7.3)	1.1 (-0.6, 2.8)	1.6 (-1.2, 4.5)
		Other Attack	5.8 (-6.2, 17.8)	2.2 (-1.9, 6.4)	3.6 (-4.3, 11.4)

Table A.2: Causal Effect Estimates and 95% Confidence Intervals for Various Stochastic Interventions. We present the results for three interventions discussed in the main text: increasing the expected number of airstrikes from 1 to 3 per day for M days, changing the focal points of airstrikes from $\alpha = 0$ to $\alpha = 3$ for M days, and the lagged effects of increasing the expected number of airstrikes from 1 to 5 per day M days ago. The range of M we consider is $\{3, 7, 30\}$. The regions of interest are Iraq, Baghdad, and the area outside Baghdad. The results in bold represent statistically significant estimates.

G.2 Empirical Results

Table A.2 presents the numerical effect estimates and 95% confidence intervals for various interventions, including those shown in the main text. We also show the effect estimates for the whole Iraq, Baghdad only, and the area outside Baghdad.

G.3 Single time point adaptive interventions

Adaptive intervention strategies are often of interest in longitudinal settings, where previous outcomes might drive future treatment assignments. In our setting, these adaptive interventions would correspond to military strategies that depend on the observed history, such as the locations of previous insurgent at-

tacks. Although we leave full development of adaptive strategies to future research, we consider adaptive strategies that take place over a single time period, and then discuss the challenges of further extending it to the multiple time period interventions.

Here, we design adaptive dosage interventions over a single time period that closely resemble the observed data in terms of the expected number of airstrikes over time and their location. Using the observed number of airstrikes over time, we fitted a smooth function of time to obtain an estimate of the expected number of airstrikes over time, which is denoted by \hat{n}_t . We used the estimated expected number of points to define adaptive interventions under which (1) the spatial distribution under the intervention is equal to the spatial distribution of airstrikes according to the propensity score, and (2) the expected number of points under the intervention is set to $c\hat{n}_t$, with c varying from 0.5 to 2 (representing a change in the number of airstrikes ranging from half to double the observed values). Formally, this intervention that depend on the observed history is given by:

$$h_{t+1}(\omega; \overline{H}_t) = \frac{c \hat{n}_t}{\int_{\Omega} h_{t+1}^{ps}(s; \overline{H}_t) ds} h_{t+1}^{ps}(\omega; \overline{H}_t),$$

where $h_{t+1}^{ps}(\omega; \overline{H}_t)$ is the estimated propensity score intensity function. This definition of intensity ensures that that expected number of airstrikes at time t is equal to $c\hat{n}_t$ using the ratio term, and the relative likelihood of each location ω being treated is as specified in the estimated propensity score. This approach is related to the incremental propensity score of [Kennedy \(2019\)](#) who considered non-spatial and non-temporal settings.

Figure A.16 shows the effect estimates for number of IED and SAF attacks in Iraq for these interventions. The result shows that the estimates are too imprecise to lead to a definitive conclusion.

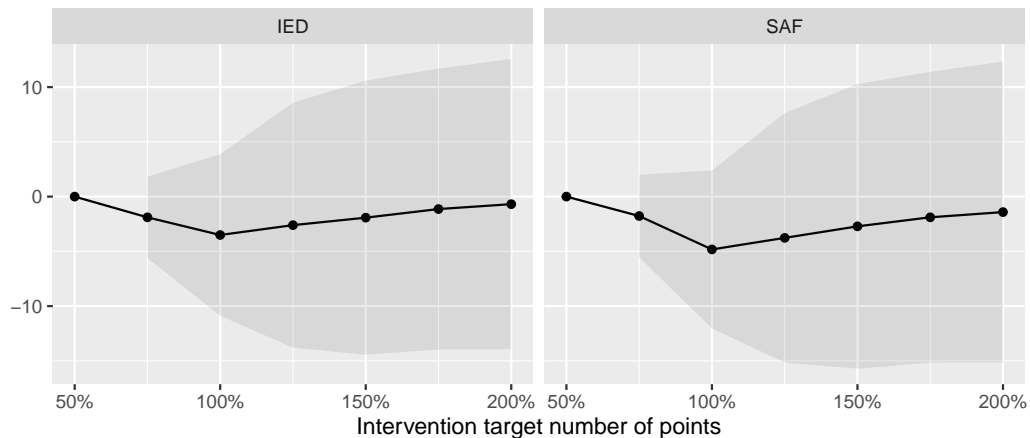


Figure A.16: Effect estimates for a change in the expected number of airstrikes from 50% \hat{n}_t to $c\hat{n}_t$, for c shown in the x-axis. Left plot shows effect estimates for IEDs and right plot shows effect estimates for SAF attacks in Iraq.

Unfortunately, the evaluation of adaptive strategies over multiple time periods rapidly becomes complicated. Specifically, for interventions over multiple time periods that depend on the most recent history, we would need to have access to intermediate potential outcomes which are unobserved. Therefore, we would have to model the outcome process in order to predict the counterfactual outcomes that would then inform the adaptive treatment assignment in the subsequent time periods. One advantage of our proposed framework is its ability to incorporate unstructured spillover and carryover effects. This is possible because our framework does not require researchers to model the outcome process. Given this difficulty, we will leave the complete investigation of adaptive spatio-temporal treatment strategies to future work.