

# Policy Evaluation for Temporal and/or Spatial Dependent Experiments

Shikai Luo<sup>a\*</sup>, Ying Yang<sup>b\*</sup>, Chengchun Shi<sup>c\*</sup>, Fang Yao<sup>d</sup>,

ridesharing applications, for instance, most AB test experiment durations do not exceed 20 days (Shi et al., 2023), and the size of treatment effects typically ranges between 0.5% and 2% (Tang et al., 2019).

The primary objective of this paper is to develop a robust statistical framework for analyzing the causal connections between the policies implemented by these companies and their corresponding outcomes, even in the presence of the aforementioned challenges. Our four major contributions can be summarized as follows. Firstly, we address the challenges by introducing linear and neural network-based Varying Coefficient Decision Process (VCDP) models. These models accommodate dynamic treatment effects over time and/or space, even in the presence of non-stationarity, random effects, interference, and spatial spillovers. These models account for market features as mediators to incorporate historical policy carryover effects. Furthermore, by assuming network interference and employing mean field approximation (as detailed in Section 3.2), we effectively operate an “effective treatment” (Manski, 2013) or “exposure mapping” (Aronow and Samii, 2017) in the spatio-temporal system. Our approach extends beyond the switchback design to any dynamic treatment allocation setup.

Secondly, we develop estimation methods for our VCDPs. For linear VCDPs, we propose a two-step process involving the calculation of least squares estimates and kernel smoothing to refine the estimates. Kernel smoothing leverages neighboring observations across time and/or space, enhancing estimation efficiency and overcoming the challenge of weak signals and small sample sizes. Additionally, we decompose average treatment effects (ATEs) into Direct Effects (DE) and Indirect Effects (IE).

et al., 2020). Our VCDPs are closely related to the second and third types of models, but they focus on interference across time *and* space. Most aforementioned works studied the interference effect across time *or* space and were motivated by research questions in environmental and epidemiological studies. It remains unknown about their generalization to ride-sharing markets. Fourthly, recent models capture the interference effect via congestion or price effects in a marketplace (Munro et al., 2021; Wager and Xu, 2021; Johari et al., 2022). These solutions rely on an assumption of Markovianity or stationarity and are design-dependent. In contrast, our approach accommodates non-stationarity and is capable of managing non-Markovianity in scenarios where outcome errors exhibit time-correlated patterns.

Our proposal is closely related to a growing literature on off-policy evaluation (OPE) methods in sequential decision making (see Uehara et al., 2022, for a review). In the literature, augmented inverse propensity score weighting methods (see e.g., Zhang et al., 2013; Luedtke and Van Der Laan, 2016; Jiang and Li, 2016; Thomas and Brunskill, 2016) have been proposed for valid OPE. Nonetheless, these methods suffer from the curse of horizon (Liu et al., 2018) in that the variance of the resulting estimator grows exponentially fast with respect to  $m$ , leading to inefficient estimates in the large  $m$  setting. Efficient model-free OPE methods have been proposed by Kallus and Uehara (2020, 2022); Liao et al. (2020, 2021); Luckett et al. (2020); Shi et al. (2021, 2022b) under the Markov decision process (MDP, see e.g., Puterman, 2014) model assumption. Recently, Hu and Wager (2021) proposed a model-free OPE method in partially observed MDPs (POMDPs) that avoids the curse of horizon. Our proposal is model-based and is ultimately different from most existing model-free OPE methods that did not consider the random effects, spatial interference effects, and the decomposition into DE and IE. In addition, little has been done on OPE for spatio-temporal dependent experiments.

Finally, our paper is related to a line of works on quantitative approaches to ride-sharing platforms. In particular, Bimpikis et al. (2019) proposed supply-and-demand models and investigated the impact of the demand pattern on the platform’s prices and profits. Castillo et al. (2017) studied how the surging prices can prevent wild goose chase (e.g., drivers pick up distant customers) and conducted regression analysis to verify the nonmonotonicity of supply on pickup times. However, estimation and inference of target policy’s treatment effect have not been considered in these papers. Cohen et al. (2022) employed the difference in differences methods to estimate the treatment effects of different types of compensation on the engagement of riders who experienced a frustration. Their analysis is limited to staggered designs. Garg and Nazerzadeh (2022) studied the theoretical properties of driver-side payment mechanisms and compared additive surge against multiplicative surge numerically. However, they did not consider the spatial spillover effects of these policies. Our paper complements the existing literature by developing a general framework to efficiently infer a target policy’s direct and indirect effects based on data collected from spatio-temporal dependent experiments and analyzing the advantage of switchback designs in the presence of spatio-temporal random effects.

## 1.2. Paper outline

The rest of the paper is organized as follows. In Section 2, we introduce a potential outcome framework for problem formulation, propose two novel temporal VCDP models under temporal dependent experiments, and develop estimation and testing procedures for both DE and IE. In Section 3, we further propose two spatio-temporal VCDP models under spatio-temporal dependent experiments and develop the associated estimation and testing procedures. In Section 4, we systematically investigate the theoretical properties of estimation and testing procedures (e.g., consistency and power) developed in Sections 2 and 3. We also illustrate the benefits of employing the switchback design in theory. In Section 5, we use numerical simulations to examine the finite sample performance of our estimation and testing procedures. Furthermore, we numerically explore the benefits of the switchback design. In Section 6, we apply the proposed procedures to evaluating different policies in Didi Chuxing.

## 2. Policy evaluation for temporal dependent experiments

In this section, we present the proposed methodology for policy evaluation in temporal dependent experiments for one experimental region.

### 2.1. A potential outcome framework

We use the potential outcome framework to present our model in non-stationary environments. We divide each day into  $m$  equally spaced nonoverlapping intervals. At each time interval, the platform can

implement either the new or old policy. We use  $A_\tau$  to denote the policy implemented at the  $\tau$ th interval for any integer  $\tau \geq 1$ . Let  $S_\tau$  be some state variables measured at the  $(\tau - 1)$ -th interval in a given day. All the states share the same support, which is assumed to be a compact subset of  $\mathbb{R}^d$ , where  $d$  denotes the dimension of the state. Let  $Y_\tau \in \mathbb{R}$  be the outcome of interest measured at time  $\tau$ .

Firstly, we define the average treatment effect (ATE) as the difference between the new and old policies. Let  $\bar{a}_\tau = (a_1, \dots, a_\tau)^\top \in \{0, 1\}^\tau$  denote a treatment history vector up to time  $\tau$ , where 1 and 0 denote the new policy and the old one, respectively. We define  $S_\tau^*(\bar{a}_{\tau-1})$  and  $Y_\tau^*(\bar{a}_\tau)$  as the counterfactual state and the counterfactual outcome, respectively. Then ATE can be defined as follows.

DEFINITION 1. *ATE is the difference between two value functions given by*

$$ATE = \sum_{\tau=1}^m \mathbb{E} [Y_\tau^*(\mathbf{1}_\tau) - Y_\tau^*(\mathbf{0}_\tau)] g,$$

where  $\mathbf{1}_\tau$  and  $\mathbf{0}_\tau$  denote vectors of 1s and 0s of length  $\tau$ , respectively.

Secondly, we can decompose ATE as the sum of direct effects (DE) and indirect effects (IE). Let  $R_\tau$  denote the conditional mean function of the outcome given the data history,

$$\mathbb{E} [Y_\tau^*(\bar{a}_\tau) | S_\tau^*(\bar{a}_{\tau-1}), Y_{\tau-1}^*(\bar{a}_{\tau-1}), S_{\tau-1}^*(\bar{a}_{\tau-2}), Y_{\tau-2}^*(\bar{a}_{\tau-2}), \dots, S_1] g = R_\tau(a_\tau, S_\tau^*(\bar{a}_{\tau-1}), a_{\tau-1}, S_\tau^*(\bar{a}_{\tau-2}), \dots, S_1).$$

LEMMA 1. Under CA, SRA and PA, we have

$$R_\tau(a_\tau, s_\tau, \dots, s_1) = \mathbb{E}(Y_\tau | A_\tau = a_\tau, S_\tau = s_\tau, \dots, S_1 = s_1), \quad (4)$$

$$\mathbb{E}fR_\tau(a, S_\tau^*(\bar{a}_{\tau-1}), \dots, S_1)g = \mathbb{E}[\mathbb{E}[R_\tau(a, S_\tau, \dots, S_1) | fA_j = a_j g_{1 \leq j < \tau}, fS_j, Y_j g_{1 \leq j < \tau}]]. \quad (5)$$

Lemma 1 implies that the causal estimand can be represented as a function of the observed data.

## 2.2. TVCDP model

We introduce two TVCDP regression models to model  $Y_{i,\tau}$  and the conditional distribution of  $S_{i,\tau}$  given the data history, forming the basis of our estimation and testing procedures. Suppose that the experiment is conducted over  $n$  days. Let  $(S_{i,\tau}, A_{i,\tau}, Y_{i,\tau})$  be the state-policy-outcome triplet measured at the  $\tau$ th time interval of the  $i$ th day for  $i = 1, \dots, n$  and  $\tau = 1, \dots, m$ . The proposed TVCDP model is composed of the following set of additive noise models,

$$\begin{aligned} Y_{i,\tau} &= f_{1,\tau}(S_{i,\tau}, A_{i,\tau}) + e_{i,\tau}, \\ S_{i,\tau+1} &= f_{2,\tau}(S_{i,\tau}, A_{i,\tau}) + \varepsilon_{i,\tau S}, \end{aligned} \quad (6)$$

where  $f_{1,\tau}(\cdot)$  and  $f_{2,\tau}(\cdot)$  are the regression functions.

We would like to highlight several key points. Firstly, in addition to defining the standard outcome regression model  $f_{1,\tau}$  as described in equation (6), it is crucial to specify how past actions influence future states. This is accomplished through the inclusion of  $f_{2,\tau}$ , which plays a pivotal role in quantifying temporal interference effects.

Secondly, we introduce a specific assumption related to the error structure. This assumption is fundamental as it allows us to incorporate temporal random effects effectively.

ASSUMPTION 1. (i) The outcome noise  $e_{i,\tau} = \eta_{i,\tau} + \varepsilon_{i,\tau}$  is a combination of two mutually independent stochastic processes: day-specific temporal variation  $\eta_{i,\tau}$  and measurement error  $\varepsilon_{i,\tau}$ . (ii) The processes  $f\eta_{i,\tau}g_{i,\tau}$  are identical realizations of a zero-mean stochastic process with covariance function  $f\Sigma_\eta(\tau_1, \tau_2)g_{\tau_1, \tau_2}$ . Additionally, all components of  $\Sigma_\eta(t_1, t_2)$  have bounded and continuous second derivatives with respect to  $t_1$  and  $t_2$ . (iii) The measurement errors  $f\varepsilon_{i,\tau}g_{i,\tau}$  and  $f\varepsilon_{i,\tau S}g_{i,\tau}$  are independent over time. They have zero mean values and exhibit  $\text{Var}(\varepsilon_{i,\tau}) = \sigma_{\varepsilon,\tau}^2$  and  $\text{Cov}(\varepsilon_{i,\tau S}) = \Sigma_{\varepsilon,\tau S}$ .

It's important to note that the day-specific random effects are present only in the outcome regression model. However, our approach can be extended to scenarios where these random effects also exist in the state regression model. We provide a detailed discussion of this extension in Section 7. Additionally, it's worth mentioning that both the conditional mean and covariance functions, namely  $f_{1,\tau}$ ,  $f_{2,\tau}$ ,  $\sigma_{\varepsilon,\tau}^2$ , and  $\Sigma_{\varepsilon,\tau S}$ , are time-dependent. This captures the nonstationarity inherent in the data generating process.

Our TVCDP models (6) have strong connections with the MDP model that is commonly used in reinforcement learning. Specifically, models (6) reduce to non-stationary (or time-varying) MDP models (Kallus and Uehara, 2022) when there are no day-specific random effects in  $f\varepsilon_{i,\tau}g_{i,\tau}$ . However, the proposed time varying models are no longer MDPs due to the existence of the day-specific random effects. In particular,  $Y_{i,\tau}$  in (6) is dependent upon past responses given  $Z_{i,\tau} = (1, S_{i,\tau}^\top, A_{i,\tau}^\top)^\top$ , leading to the violation of the conditional independence assumption. In addition, the market features at each time serve as mediators that mediate the effects of past actions on the current outcome.

Next, we consider two specific function approximations for  $f_1$  and  $f_2$  and derive their related IE and DE as follows.

MODEL 1. Linear temporal varying coefficient decision process (L-TVCDP) assumes

$$\begin{aligned} Y_{i,\tau} &= \beta_0(\tau) + S_{i,\tau}^\top \beta(\tau) + A_{i,\tau} \gamma(\tau) + e_{i,\tau} = Z_{i,\tau}^\top \theta(\tau) + e_{i,\tau}, \\ S_{i,\tau+1} &= \phi_0(\tau) + \Phi(\tau) S_{i,\tau} + A_{i,\tau} \Gamma(\tau) + \varepsilon_{i,\tau S} = \Theta(\tau) Z_{i,\tau} + \varepsilon_{i,\tau S}, \end{aligned}$$

where  $\theta(\tau) = (\beta_0(\tau), \beta(\tau)^\top, \gamma(\tau))^\top$  is a  $(d+2) \times 1$  vector of time-varying coefficients,  $\Theta(\tau) = [\phi_0(\tau) \ \Phi(\tau) \ \Gamma(\tau)]$  is a  $d \times (d+2)$  coefficient matrix and  $Z_{i,\tau} = (1, S_{i,\tau}^\top, A_{i,\tau}^\top)^\top$ .

Model 1 shares a close connection with the linear quadratic Gaussian model (LQG), well studied in the fields of RL and control theory (see, for example, Lale et al., 2021). To be more specific, Model 1 can be seen as a simplified, one-dimensional observation variant of LQG under certain conditions. This happens when the outcome regression model doesn't incorporate  $A_{i,\tau}$  and the autocorrelated noise  $\eta_{i,\tau}$ . However,

there's a crucial distinction between LQG and our proposed model. In LQG, the state variables are hidden and must be deduced from the observed  $Y_{i,\tau}$  values. This contrasts with similar models used in literature for estimating dynamic treatment effects (Lewis and Syrgkanis, 2020).

When  $f\eta_{i,\tau}g_{i,\tau}$  become the fixed effects and satisfy  $\eta_{i,\tau} = \eta_i$  for any  $i$  and  $\tau$ , the outcome regression model of L-TVCDP includes both the day-specific fixed effects  $f\eta_i g_i$  and the time-specific fixed effects  $f\beta_0(\tau)g_\tau$ . It is similar to the two-way fixed effects model in the panel data literature (De Chaisemartin and d'Haultfoeuille, 2020; Wooldridge, 2021; Arkhangelsky et al., 2021; Imai and Kim, 2021). Furthermore, we derive the closed-form expressions for DE and IE under L-TVCDP, whose proof can be found in Section S.3 of the supplementary document.

**PROPOSITION 1.** *Under the L-TVCDP model, we have  $DE = \sum_{\tau=1}^m \gamma(\tau)$  and*

$$IE = \sum_{\tau=2}^m \beta(\tau)^\top \left\{ \sum_{k=1}^{\tau-1} (\Phi(\tau-1)\Phi(\tau-2)\dots\Phi(k+1)) \Gamma(k) \right\}, \quad (7)$$

where by convention, the product  $\Phi(\tau-1)\Phi(\tau-2)\dots\Phi(k+1) = 1$  when  $\tau-1 < k+1$ .

**MODEL 2.** *Neural networks temporal varying decision process (NN-TVCDP) assumes*

$$\begin{aligned} Y_{i,\tau} &= g_0(\tau, S_{i,\tau}) \mathbb{I}(A_{i,\tau} = 0) + g_1(\tau, S_{i,\tau}) \mathbb{I}(A_{i,\tau} = 1) + e_{i,\tau}, \\ S_{i,\tau+1} &= G_0(\tau, S_{i,\tau}) \mathbb{I}(A_{i,\tau} = 0) + G_1(\tau, S_{i,\tau}) \mathbb{I}(A_{i,\tau} = 1) + \varepsilon_{i,\tau} S, \end{aligned}$$

where  $\mathbb{I}(\cdot)$  denotes the indicator function of an event and  $g_0(\cdot, \cdot)$ ,  $g_1(\cdot, \cdot)$ ,  $G_0(\cdot, \cdot)$ , and  $G_1(\cdot, \cdot)$  are parametrized via some (deep) neural networks.

Under NN-TVCDP, DE and IE are, respectively, given by

$$DE = \sum_{\tau=1}^m \mathbb{E} \{ g_1(\tau, S_\tau^0) - g_0(\tau, S_\tau^0) \} \quad \text{and} \quad IE = \sum_{\tau=1}^m \mathbb{E} \{ g_1(\tau, S_\tau^1) - g_1(\tau, S_\tau^0) \}, \quad (8)$$

where  $S_\tau^0$  and  $S_\tau^1$  are defined recursively by  $S_\tau^0 = G_0(\tau-1, S_{\tau-1}^0)$  and  $S_\tau^1 = G_1(\tau-1, S_{\tau-1}^1)$ .

### 2.3. Estimation and testing procedures for DE in the L-TVCDP model

We describe our estimation and testing procedures for DE in the L-TVCDP model and present their pseudocode in Algorithm 1 as follows.

---

#### **Algorithm 1** Inference of DE in the L-TVCDP model

---

- 1: Compute the OLS estimator  $\hat{\theta}$  according to (9).
  - 2: Employ kernel smoothing to compute a refined estimator  $\tilde{\theta}$  according to (10) and calculate the estimate  $\widehat{DE}$  by (11).
  - 3: Estimate the variance of  $\hat{\theta}$  as follows:
  - 4: (3.1). Estimate the conditional variance of  $\mathbf{Y}_i$  given  $fZ_{i,\tau}g_\tau$  using (12);
  - 5: (3.2). Estimate the variance of  $\hat{\theta}$  by the sandwich estimator (13).
  - 6: Estimate the variance of  $\tilde{\theta}$  by  $\tilde{\mathbf{V}}_\theta = \mathbf{\Omega} \widehat{\mathbf{V}}_\theta \mathbf{\Omega}^\top$  and compute the standard error of  $\widehat{DE}$ , denoted by  $\widehat{se}(\widehat{DE})$ .
  - 7: Reject  $H_0^{DE}$  if  $\widehat{DE}/\widehat{se}(\widehat{DE})$  exceeds the upper  $\alpha$ th quantile of a standard normal distribution.
- 

Step 1 of Algorithm 1 is to obtain an initial estimator of  $\theta(\tau)$  by computing its ordinary least squares (OLS) estimator, defined as

$$\hat{\theta}(\tau) = \left( \sum_{i=1}^n Z_{i,\tau} Z_{i,\tau}^\top \right)^{-1} \left( \sum_{i=1}^n Z_{i,\tau} Y_{i,\tau} \right) \quad \text{for } 1 \leq \tau \leq m. \quad (9)$$

Step 2 of Algorithm 1 is to employ kernel smoothing to refine the initial estimator. Specifically, for a given kernel function  $K(\cdot)$ , we introduce the refined estimator

$$\tilde{\theta}(\tau) = (\tilde{\beta}_0(\tau), \tilde{\beta}(\tau)^\top, \tilde{\gamma}(\tau)^\top)^\top = \sum_{\tau=1}^m \omega_{\tau,h}(t) \hat{\theta}(\tau), \quad (10)$$

for any  $t \in [0, m]$  and a bandwidth parameter  $h$ , where  $\omega_{\tau,h}(t) = K((t - \tau)/(mh)) / \sum_{j=1}^m K((t - j)/(mh))$  is the weight function. Our DE estimator is given by

$$\widehat{\text{DE}} = \sum_{\tau=1}^m \tilde{\gamma}(\tau). \quad (11)$$

We will show in Section 4 that as  $\min(n, m) \rightarrow \infty$ ,  $\widehat{\text{DE}}$  is asymptotically normal. To derive a Wald test for (2), it remains to estimate its variance  $\text{Var}(\widehat{\text{DE}})$ .

There are two major advantages of using the smoothing step here. First, it allows us to estimate the time-varying coefficient curve  $\theta(t)$  without restricting  $t$  to the class of integers. Second, the smoothed estimator has smaller variance, leading to a more powerful test statistics. To elaborate, according to model (6) for L-TVCDP, the variation of the OLS estimator comes from two sources, the day-specific random effect and the measurement error. The use of smoothing removes the random fluctuations due to the measurement error. See Theorem 1 in Section 4 for a formal statement. This smoothing technique has been widely applied in the analysis of varying-coefficient models (see e.g., Zhu et al., 2014).

Step 3 of Algorithm 1 is to estimate the covariance matrix of the initial estimator  $\hat{\theta} = (\hat{\theta}^\top(1), \dots, \hat{\theta}^\top(m))^\top$ . We first estimate the residual  $e_{i,\tau}$  by  $\hat{e}_{i,\tau} = Y_{i,\tau} - \mathbf{Z}_{i,\tau}^\top \tilde{\theta}(\tau)$ . It allows us to estimate the day-specific random effect via smoothing, i.e.,  $\hat{\eta}_i(t) = \sum_{j=1}^m \omega_{j,h}(t) \hat{e}_{i,\tau}$ . Second, the measurement error can be estimated by  $\hat{\varepsilon}_{i,\tau} = \hat{e}_{i,\tau} - \hat{\eta}_{i,\tau}$  for any  $i$  and  $\tau$ , where  $\hat{\eta}_{i,\tau} = \hat{\eta}_i(\tau)$ . Third, we estimate the conditional covariance matrix of  $\mathbf{Y}_i = (Y_{i,1}, \dots, Y_{i,m})^\top$  given  $f\mathbf{Z}_{i,\tau}g_\tau$  based on these estimated residuals. Under model (6) for L-TVCDP, the covariance between  $Y_{i,\tau_1}$  and  $Y_{i,\tau_2}$  conditional on  $f\mathbf{Z}_{i,\tau}g_\tau$  is given by  $\Sigma_y(\tau_1, \tau_2) = \sigma_{\varepsilon,\tau_1}^2 \mathbb{I}(\tau_1 = \tau_2) + \Sigma_\eta(\tau_1, \tau_2)$ , which can be consistently estimated by

$$\hat{\Sigma}_y(\tau_1, \tau_2) = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_{i,\tau_1}^2 \mathbb{I}(\tau_1 = \tau_2) + \frac{1}{n} \sum_{i=1}^n \hat{\eta}_{i,\tau_1} \hat{\eta}_{i,\tau_2}. \quad (12)$$

This allows us to estimate  $\text{Var}(\mathbf{Y}_i | f\mathbf{Z}_{i,\tau}g_\tau)$  by  $\hat{\Sigma} = f\hat{\Sigma}_y(\tau_1, \tau_2)g_{\tau_1, \tau_2}$ . Finally, the covariance matrix of  $\hat{\theta}$  can be consistently estimated by the sandwich estimator,

$$\hat{\mathbf{V}}_\theta = \left( \sum_{i=1}^n \mathbf{Z}_i^\top \mathbf{Z}_i \right)^{-1} \left( \sum_{i=1}^n \mathbf{Z}_i^\top \hat{\Sigma} \mathbf{Z}_i \right) \left( \sum_{i=1}^n \mathbf{Z}_i^\top \mathbf{Z}_i \right)^{-1}, \quad (13)$$

where  $\mathbf{Z}_i$  is a block-diagonal matrix computed by aligning  $\mathbf{Z}_{i,1}^\top, \dots, \mathbf{Z}_{i,m}^\top$  along its diagonal.

Step 4 of Algorithm 1 is to estimate the covariance matrix of the refined estimator  $\tilde{\theta} = (\tilde{\theta}^\top(1), \dots, \tilde{\theta}^\top(m))^\top$ . A key observation is that each  $\tilde{\theta}(\tau)$  is essentially a weighted average of  $f\hat{\theta}(\tau)g_\tau$ . Writing in matrix form, we have  $\tilde{\theta} = \Omega \hat{\theta}$ , where  $\Omega$  is a block-diagonal matrix computed by aligning  $\omega_{1,h}(\tau) \mathbf{J}_p, \dots, \omega_{m,h}(\tau) \mathbf{J}_p$  along its diagonal and  $\mathbf{J}_p$  is a  $p \times p$  matrix of ones. As such, we estimate the covariance matrix of  $\tilde{\theta}$  by  $\hat{\mathbf{V}}_{\tilde{\theta}} = \Omega \hat{\mathbf{V}}_\theta \Omega^\top$ . This in turn yields a consistent estimator for the variance of  $\widehat{\text{DE}}$ , as  $\widehat{\text{DE}}$  is a linear combination of  $\tilde{\theta}$ .

Step 5 of Algorithm 1 is to construct a Wald-type test statistic based on  $\widehat{\text{DE}}$  and its standard error  $\widehat{\text{se}}(\widehat{\text{DE}})$ . We reject the null hypothesis in (2) if  $\widehat{\text{DE}}/\widehat{\text{se}}(\widehat{\text{DE}})$  exceeds the upper  $\alpha$ th quantile of a standard normal distribution. Size and power properties of the proposed test are investigated in Section 4.

#### 2.4. Estimation and testing procedures for IE in the L-TVCDP model

We describe our estimation and testing procedures for IE in the L-TVCDP model and present their pseudocode in Algorithm 2 as follows.

---

##### Algorithm 2 Inference of IE in the L-TVCDP model

---

- 1: Compute the OLS estimator

$$\hat{\Theta} = f\hat{\Theta}(1), \dots, \hat{\Theta}(m-1)g^\top = f \sum_{i=1}^n \mathbf{Z}_{i,(-m)} \mathbf{Z}_{i,(-m)}^\top g^{-1} f \sum_{i=1}^n \mathbf{Z}_{i,(-m)} \mathbf{S}_{i,(-1)}^\top g,$$

where  $\mathbf{S}_{i,(-1)}$  and  $\mathbf{Z}_{i,(-m)}$  are block-diagonal matrices computed by aligning  $\mathbf{S}_{i,2}^\top, \dots, \mathbf{S}_{i,m}^\top$  and  $\mathbf{Z}_{i,1}^\top, \dots, \mathbf{Z}_{i,m-1}^\top$  along their diagonals, respectively.

- 2: Compute the refined estimator  $\tilde{\Theta} = f\tilde{\Theta}(1), \dots, \tilde{\Theta}(m-1)g^\top = \Omega\hat{\Theta}$ .
- 3: Construct the plug-in estimator  $\widehat{\text{IE}}$  according to (14).
- 4: Compute the estimated residual  $\hat{\varepsilon}_{i,\tau S} = S_{i,\tau+1} - Z_{i,\tau}\tilde{\Theta}(\tau)$  for any  $i$  and  $\tau$ .
- 5: **for**  $b = 1, \dots, B$  **do**  
 Generate i.i.d. standard normal random variables  $f\xi_i^b g_{i=1}^n$ ;  
 Generate pseudo outcomes  $f\hat{S}_{i,\tau}^b g_{i,\tau}$  and  $f\hat{Y}_{i,\tau}^b g_{i,\tau}$  according to (15);  
 Repeat Steps 1-2 in Algorithm 1 and Steps 1-3 in Algorithm 2 to compute  $\widehat{\text{IE}}^b$ .
- 6: **end for**
- 7: Reject  $H_0^{IE}$  if  $\widehat{\text{IE}}$  exceeds the upper  $\alpha$ th empirical quantile of  $f\widehat{\text{IE}}^b - \widehat{\text{IE}}g_b$ .

Steps 1-3 of Algorithm 2 are to compute a consistent estimator  $\widehat{\text{IE}}$  for IE. Specifically, in Step 1 of Algorithm 2, we apply OLS regression to derive an initial estimator  $\hat{\Theta}$  for  $\Theta = f\Theta(1), \dots, \Theta(m-1)g^\top$ . In Step 2 of Algorithm 2, we employ kernel smoothing to compute a refined estimator  $\tilde{\Theta} = \Omega\hat{\Theta}$  to improve its statistical efficiency, as in Algorithm 1. In Step 3 of Algorithm 2, we plug in  $\tilde{\Theta}$  and  $\tilde{\theta}$  for  $\Theta$  and  $\theta$  in model 1, leading to

$$\widehat{\text{IE}} = \sum_{\tau=2}^m \tilde{\beta}(\tau)^\top \left\{ \sum_{k=1}^{\tau-1} \left( \tilde{\Phi}(\tau-1)\tilde{\Phi}(\tau-2) \dots \tilde{\Phi}(k+1) \right) \tilde{\Gamma}(k) \right\}, \quad (14)$$

where  $\tilde{\beta}(\tau)$ ,  $\tilde{\Phi}(\tau)$  and  $\tilde{\Gamma}(\tau)$  are the corresponding estimators for  $\beta(\tau)$ ,  $\Phi(\tau)$  and  $\Gamma(\tau)$ , respectively.

Step 4 of Algorithm 2 is to compute the estimated residuals  $\hat{E}_{i,\tau} = S_{i,\tau+1} - Z_{i,\tau}\tilde{\Theta}(\tau)$  for all  $i$  and  $\tau$ , which are used to generate pseudo outcomes in the subsequent bootstrap step.

Step 5 of Algorithm 2 is to use bootstrap to simulate the distribution of  $\widehat{\text{IE}}$  under the null hypothesis. The key idea is to compute the bootstrap samples for  $\tilde{\theta}$  and  $\tilde{\Theta}$  and use the plug-in principle to construct the bootstrap samples for  $\widehat{\text{IE}}$ . A key observation is that  $\tilde{\theta}$  and  $\tilde{\Theta}$  depend linearly on the random errors, so the wild bootstrap method (Wu et al., 1986) is applicable. We begin by generating i.i.d. standard normal random variables  $f\xi_i g_{i=1}^n$ . We next generate pseudo-outcomes given by

$$\hat{S}_{i,\tau+1} = \tilde{\Theta}(\tau)\hat{Z}_{i,\tau} + \xi_i\hat{\varepsilon}_{i,\tau S} \text{ and } \hat{Y}_{i,\tau} = \hat{Z}_{i,\tau}^\top\tilde{\theta}(\tau) + \xi_i\hat{\varepsilon}_{i,\tau}, \quad (15)$$

where  $\hat{Z}_{i,\tau}$  is a version of  $Z_{i,\tau}$  with  $S_{i,\tau}$  replaced by  $\hat{S}_{i,\tau}$ . Furthermore, we apply Steps 1-2 of Algorithm 1 and Steps 1-3 of Algorithm 2 to compute the bootstrap version of  $\widehat{\text{IE}}$  based on these pseudo outcomes in (15). The above procedures are repeatedly applied to simulate a sequence of bootstrap estimators  $f\widehat{\text{IE}}^b g_{b=1}^B$  based on which the decision region can be derived.

### 2.5. Estimation procedure in NN-TVCDP model

We first introduce how to estimate the regression functions  $g_0, g_1, G_0$  and  $G_1$ . Take  $g_0$  as an instance, we consider minimizing the following empirical objective function

$$\sum_{i=1}^n \sum_{\tau=1}^m (1 - A_{i,\tau}) fY_{i,\tau} - g_0(\tau, S_{i,\tau})g^2.$$

Instead of separately estimating  $g_0(\tau, \cdot)$  for each  $\tau$ , we treat  $\tau$  as part of the features and jointly estimate  $f\hat{g}_0(\tau, \cdot)g_\tau$  by solving the above optimization. It allows us to borrow information across different time points to improve the estimation accuracy.

Next, we introduce the estimation procedures for DE and IE. We impose a parametric model (e.g., Gaussian) for the density function  $f_{\varepsilon_{\tau S}}$  of the measurement error  $\varepsilon_{i,\tau S}$  and summarize the steps below.

1. Use neural networks to estimate  $g_0, g_1, G_0$  and  $G_1$  by solving their corresponding least square objective functions. Denote the corresponding estimators by  $\hat{g}_0, \hat{g}_1, \hat{G}_0$ , and  $\hat{G}_1$ , respectively.
2. Compute the residual  $\hat{\varepsilon}_{i,\tau S} = S_{i,\tau+1} - \left\{ \hat{G}_0(\tau, S_{i,\tau}) \mathbb{I}(A_{i,\tau} = 0) + \hat{G}_1(\tau, S_{i,\tau}) \mathbb{I}(A_{i,\tau} = 1) \right\}$  and use  $\hat{\varepsilon}_{i,\tau S}$  to compute the density function estimator  $\hat{f}_{\varepsilon_{\tau S}}$ .



3. Use Monte Carlo to estimate the distributions of the potential states  $S_{i,\tau}^*(\mathbf{1}_{\tau-1})$  and  $S_{i,\tau}^*(\mathbf{0}_{\tau-1})$  conditional on  $S_{i,1}$ . Specifically, for  $\tau = 1, \dots, m$ ,  $i = 1, \dots, n$ , and  $k = 1, \dots, M$ , we use  $\widehat{f}_{\varepsilon_{\tau S}}$  to generate error residuals  $\widehat{f}_{\widehat{\varepsilon}_{i,\tau S,k}} \mathcal{G}_{k=1}^M$ , where  $M$  denotes the number of Monte Carlo replications. Next, we set  $\widehat{S}_{i,1,k}^1 = \widehat{S}_{i,1,k}^0 = S_{i,1}$  for any  $i$  and  $k$ , and sequentially construct Monte Carlo samples  $f_{\widehat{S}_{i,\tau,k}^1} \mathcal{G}_{k=1}^M, f_{\widehat{S}_{i,\tau,k}^0} \mathcal{G}_{k=1}^M$  by setting  $\widehat{S}_{i,\tau+1,k}^1 = \widehat{G}_1(\tau, \widehat{S}_{i,\tau,k}^1) + \widehat{\varepsilon}_{i,\tau S,k}$  and  $\widehat{S}_{i,\tau+1,k}^0 = \widehat{G}_0(\tau, \widehat{S}_{i,\tau,k}^0) + \widehat{\varepsilon}_{i,\tau S,k}$ .
4. Based on (8), we estimate DE and IE by using

$$\begin{aligned} \widehat{DE} &= \frac{1}{nM} \sum_{i=1}^n \sum_{k=1}^M \sum_{\tau=1}^m \left\{ \widehat{g}_1(\tau, \widehat{S}_{i,k,\tau}^0) \quad \widehat{g}_0(\tau, \widehat{S}_{i,k,\tau}^0) \right\} \quad \text{and} \\ \widehat{IE} &= \frac{1}{nM} \sum_{i=1}^n \sum_{k=1}^M \sum_{\tau=2}^m \left\{ \widehat{g}_1(\tau, \widehat{S}_{i,k,\tau}^1) \quad \widehat{g}_1(\tau, \widehat{S}_{i,k,\tau}^0) \right\}. \end{aligned}$$

### 3. Policy evaluation for spatio-temporal dependent experiments

In this section, we present the proposed methodology for policy evaluation in spatio-temporal dependent experiments by extending our proposal in temporal dependent experiments. We highlight several key differences between the spatio-temporal dependent experiment and the temporal dependent one.

#### 3.1. A potential outcome framework

Firstly, we introduce the spatio-temporal dependent experiments as follows. Specifically, a city is split into  $r$  non-overlapping regions. Each region receives a sequence of policies over time and different regions may receive different policies at the same time. In our application, we employ the spatio-temporal dependent alternation design to randomize these policies. In each region, we independently randomize the initial policy (either A or B) and then apply the temporal alternation design. As discussed in the introduction, one major challenge for policy evaluation is that the spatial proximities will induce spatio-temporal interference among locations across time. In the example of ride-sharing platforms, for many call orders, their pickup locations and destinations belong to different regions. Therefore, applying an order dispatch policy at one region will change the distribution of drivers of its neighbouring areas as well, so the order dispatch policy at one location could influence outcomes of those neighbouring areas, inducing interference among spatial units.

Secondly, to quantify the spatio-temporal interference, we allow the potential outcome of each region to depend on policies applied to its neighbouring areas as well. Specifically, for the  $\iota$ th region, let  $\bar{a}_{\tau,\iota} = (\bar{a}_{1,\iota}, \dots, \bar{a}_{\tau,\iota})^\top$  denote its treatment history up to time  $\tau$  and  $N_\iota$  denote the neighbouring regions of  $\iota$ . Let  $\bar{a}_{\tau,[1:r]} = (\bar{a}_{\tau,1}, \dots, \bar{a}_{\tau,r})^\top$  denote the treatment history associated with all regions. Similarly, let  $S_{\tau,\iota}^*(\bar{a}_{\tau-1,[1:r]})$  and  $Y_{\tau,[1:r]}^*(\bar{a}_{\tau,[1:r]})$  denote the potential state and outcome associated with the  $\iota$ th region, respectively. Let  $S_{\tau,[1:r]}^*(\bar{a}_{\tau-1,[1:r]})$  denote the set of potential states at time  $\tau$ .

Similarly, we introduce CA and SRA in the spatio-temporal case as follows.

**CA.**  $S_{\tau+1,\iota}^*(\bar{A}_{\tau,[1:r]}) = S_{\tau+1,\iota}$  and  $Y_{\tau,\iota}^*(\bar{A}_{\tau,[1:r]}) = Y_{\tau,\iota}$  for any  $\tau \geq 1$  and  $1 \leq \iota \leq r$ , where  $\bar{A}_{\tau,[1:r]}$  denotes the set of observed treatment history up to time  $\tau$ .

**SRA.**  $A_{\tau,[1:r]}$ , the set of observed policies at time  $\tau$ , is conditionally independent of all potential variables given  $S_{\tau,[1:r]}$  and  $f(S_{j,[1:r]}, A_{j,[1:r]}, Y_{j,[1:r]}) \mathcal{G}_{j < \tau}$ .

SRA automatically holds under the spatio-temporal alternation design, in which the policy assignment mechanism is conditionally independent of the data given the policies assigned at the initial time point.

Thirdly, we are interested in the overall treatment effects. Define ATE as the difference between the new and old policies aggregated over different regions.

DEFINITION 2. ATE is defined as the difference between two value functions given by

$$ATE_{st} = \sum_{\iota=1}^r \sum_{\tau=1}^m \mathbb{E} f Y_{\tau,\iota}^*(\mathbf{1}_{\tau,[1:r]}) - Y_{\tau,\iota}^*(\mathbf{0}_{\tau,[1:r]}) \mathcal{G}_{\tau}.$$

Let  $R_{\tau,\ell}$  denote the conditional mean function of  $Y_{\tau,\ell}^*(\bar{a}_{\tau,[1:r]})$  given the past policies and potential states. Similarly, we can decompose ATE as the sum of DE and IE, which are, respectively, given by

$$\begin{aligned} \text{DE}_{st} &= \sum_{\ell=1}^r \sum_{\tau=1}^m \mathbb{E} f R_{\tau,\ell}(\mathbf{1}_{\tau,[1:r]}, S_{\tau,\ell}^*(\mathbf{0}_{\tau-1,[1:r]}), \mathbf{0}_{\tau-1,[1:r]}, \dots, S_1) R_{\tau,\ell}(\mathbf{0}_{\tau,[1:r]}, S_{\tau,\ell}^*(\mathbf{0}_{\tau-1,[1:r]}), \mathbf{0}_{\tau-1,[1:r]}, \dots, S_1) g, \\ \text{IE}_{st} &= \sum_{\ell=1}^r \sum_{\tau=1}^m \mathbb{E} f R_{\tau,\ell}(\mathbf{1}_{\tau,[1:r]}, S_{\tau,\ell}^*(\mathbf{1}_{\tau-1,[1:r]}), \mathbf{1}_{\tau-1,[1:r]}, \dots, S_1) R_{\tau,\ell}(\mathbf{1}_{\tau,[1:r]}, S_{\tau,\ell}^*(\mathbf{0}_{\tau-1,[1:r]}), \mathbf{0}_{\tau-1,[1:r]}, \dots, S_1) g. \end{aligned}$$

We aim to test the following hypotheses:

$$H_0^{DE} : \text{DE}_{st} = 0 \text{ v.s. } H_1^{DE} : \text{DE}_{st} > 0, \quad (16)$$

$$H_0^{IE} : \text{IE}_{st} = 0 \text{ v.s. } H_1^{IE} : \text{IE}_{st} > 0. \quad (17)$$

### 3.2. Spatio-temporal VCDP models

We introduce the spatio-temporal VCDP (STVCDP) models to model  $Y_{\tau,\ell}$  and  $S_{\tau,\ell}$ , respectively. Suppose that the experiment is conducted across  $r$  regions over  $n$  days. Let  $(S_{i,\tau,\ell}, A_{i,\tau,\ell}, Y_{i,\tau,\ell})$  denote the state-policy-outcome triplet measured from the  $\ell$ th region at the  $\tau$ th time interval of the  $i$ th day for  $i = 1, \dots, n$ ,  $\tau = 1, \dots, m$ , and  $\ell = 1, \dots, r$ . The STVCDP model is given as follows,

$$\begin{aligned} Y_{i,\tau,\ell} &= f_{1,\tau,\ell}(S_{i,\tau,\ell}, A_{i,\tau,\ell}, \bar{A}_{i,\tau,\ell}) + e_{i,\tau,\ell}, \\ S_{i,\tau+1,\ell} &= f_{2,\tau,\ell}(S_{i,\tau,\ell}, A_{i,\tau,\ell}, \bar{A}_{i,\tau,\ell}) + \epsilon_{i,\tau,\ell}, \end{aligned}$$

where  $\bar{A}_{i,\tau,\ell}$  denotes the average of  $fA_{i,\tau,k}g_{k \in \mathcal{N}_\ell}$ , and  $f e_{i,\tau,\ell}, \epsilon_{i,\tau,\ell}g$  are the random noises. In parallel to Assumption 1, we impose the following noise assumption for the STVCDP model.

**ASSUMPTION 2.** (i) The outcome noise  $e_{i,\tau,\ell} = \eta_{i,\tau,\ell}^I + \eta_{i,\tau,\ell}^{II} + \eta_{i,\tau,\ell}^{III} + \varepsilon_{i,\tau,\ell}$  can be decomposed into four mutually independent processes:  $f\eta_{i,\tau,\ell}^I g, f\eta_{i,\tau,\ell}^{II} g, f\eta_{i,\tau,\ell}^{III} g$ , and  $f\varepsilon_{i,\tau,\ell} g$ . (ii) The  $f\eta_{i,\tau,\ell}^I g, f\eta_{i,\tau,\ell}^{II} g$  and  $f\eta_{i,\tau,\ell}^{III} g$  are i.i.d. copies of some zero-mean random processes with covariance functions  $\Sigma_{\eta^I}(\tau_1, \ell_1, \tau_2, \ell_2), \Sigma_{\eta^{II}}(\tau_1, \ell_1, \tau_2, \ell_2)\mathbb{I}(\ell_1 = \ell_2)$ , and  $\Sigma_{\eta^{III}}(\tau_1, \ell_1, \tau_2, \ell_2)\mathbb{I}(\tau_1 = \tau_2)$ , respectively. These covariance functions have bounded and continuously differentiable second-order derivatives. (iii) The measurement errors  $f\varepsilon_{i,\tau,\ell}g_{i,\tau,\ell}$  and the state noises  $f\epsilon_{i,\tau,\ell}g_{i,\tau,\ell}$  are independent over different location/time combinations, have zero means, and satisfy  $\text{Var}(\varepsilon_{i,\tau,\ell}) = \sigma_\varepsilon^2(\tau, \ell)$  and  $\text{Cov}(\epsilon_{i,\tau,\ell}) = \Sigma_{\epsilon,\tau,\ell}$ .

We make three remarks. Firstly, as per the STVCDP model, the outcome in the  $\ell$ th region is influenced solely by the current actions  $A_{i,\tau,\ell}$  and those from its neighboring areas. This assumption is often valid in various applications, such as ride-sharing platforms. For instance, the policy in one location may impact other locations only through its effect on the distribution of drivers. Within each time unit, a driver can travel at most from one location to its neighboring ones. Consequently, outcomes in one location are independent of policies applied to non-adjacent locations.

Secondly, in our spatial interference model, we adopt the mean field approximation. Under this approach, the outcome  $Y_{\tau,\ell}$  and next state  $S_{\tau+1,\ell}$  in a given region depend on the treatments of neighboring regions  $fA_{\tau,k}g_{k \in \mathcal{N}_\ell}$  only through their average  $\bar{A}_{\tau,\ell}$ . The mean field approximation is a commonly used technique in multi-agent reinforcement learning for policy learning and evaluation. It's worth noting that studies, such as Shi et al. (2022a), have shown that the average effect  $\bar{A}_{\tau,\ell}$  effectively summarizes the impact of  $fA_{\tau,k}g_{k \in \mathcal{N}_\ell}$ . This approach aligns with assumptions frequently made in the causal inference literature dealing with spatial interference (Sobel, 2006; Hudgens and Halloran, 2008; Zigler et al., 2012; Perez-Heydrich et al., 2014; Sobel and Lindquist, 2014; Liu et al., 2016; Sävje et al., 2021).

Thirdly, besides the average effect, alternative low-dimensional summary statistics of  $fA_{ij} : j \in \mathcal{N}_\ell g$  can be considered, such as  $\sum_{j \in \mathcal{N}_\ell} \theta_{ij} A_{ij}$  and  $\theta_\ell \mathbb{I}_{\{\sum_{j \in \mathcal{N}_\ell} A_{ij} > 0\}}$  (Hu et al., 2022). The resulting estimation and inference procedures can be similarly derived.

Similar to model (6), we allow general function approximation for  $f_1$  and  $f_2$ . To save space, we focus on linear STVCDP models (L-STVCDP) in the rest of this section. Meanwhile, the proposed estimation procedure can be extended to handle neural network STVCDP models, as in Section 2.4. The proposed L-STVCDP model is given as follows,

$$\begin{aligned} Y_{i,\tau,\ell} &= \beta_0(\tau, \ell) + S_{i,\tau,\ell}^\top \beta(\tau, \ell) + A_{i,\tau,\ell} \gamma_1(\tau, \ell) + \bar{A}_{i,\tau,\ell} \gamma_2(\tau, \ell) + e_{i,\tau,\ell}, \\ S_{i,\tau+1,\ell} &= \phi_0(\tau, \ell) + \Phi(\tau, \ell) S_{i,\tau,\ell} + A_{i,\tau,\ell} \Gamma_1(\tau, \ell) + \bar{A}_{i,\tau,\ell} \Gamma_2(\tau, \ell) + \epsilon_{i,\tau,\ell}, \end{aligned} \quad (18)$$

where  $Z_{i,\tau,\iota} = (1, S_{i,\tau,\iota}^\top, A_{i,\tau,\iota}, \bar{A}_{i,\tau,\mathcal{N}_i})^\top$ .

Similar to (7), we can show that  $\text{DE}_{st}$  and  $\text{IE}_{st}$  are equal to the following,

$$\begin{aligned} \text{DE}_{st} &= \sum_{\iota=1}^r \sum_{\tau=1}^m \tilde{f}\gamma_1(\tau, \iota) + \gamma_2(\tau, \iota)g, \\ \text{IE}_{st} &= \sum_{\iota=1}^r \sum_{\tau=1}^m \beta(\tau, \iota)^\top \left[ \sum_{k=1}^{\tau-1} (\Phi(\tau-1, \iota) \dots \Phi(k+1, \iota)) \tilde{f}\Gamma_1(k, \iota) + \Gamma_2(k, \iota)g \right], \end{aligned} \quad (19)$$

where the product  $\Phi(\tau-1, \iota) \dots \Phi(k+1, \iota) = 1$  when  $\tau-1 < k+1$ . These two identities form the basis of our test procedure.

### 3.3. Estimation and testing procedures for DE and IE

We first describe our estimation and testing procedures for DE under the spatio-temporal alternation design and present the pseudocode in Algorithms S.1 of Section S.1 of the supplementary document to save space.

Step 1 of Algorithm S.1 is to independently apply Steps 1 and 2 of Algorithm 1 detailed in Section 2.3 to the data subset  $\tilde{f}(Z_{i,\tau,\iota}, Y_{i,\tau,\iota})g_{i,\tau}$  for each region  $\iota$  in order to compute a smoothed estimator  $\tilde{\theta}_{st}^0(\iota) = \tilde{f}\tilde{\theta}_{st}^0(1, \iota)^\top, \dots, \tilde{\theta}_{st}^0(m, \iota)^\top g^\top$  for  $\tilde{f}\theta(1, \iota)^\top, \dots, \theta(m, \iota)^\top g^\top$ .

Step 2 of Algorithm S.1 is to employ kernel smoothing again to spatially smooth each component of  $\tilde{\theta}_{st}^0(\iota)$  across all  $\iota \in \{1, \dots, rg\}$ . Specifically, we compute  $\tilde{\theta}_{st}(\iota) = \tilde{f}\tilde{\theta}_{st}(1, \iota)^\top, \dots, \tilde{\theta}_{st}(m, \iota)^\top g^\top$  as the resulting refined estimator, given by  $\tilde{\theta}_{st}(\tau, \iota) = \sum_{\ell=1}^r \kappa_{\ell, h_{st}}(\iota) \tilde{\theta}_{st}^0(\tau, \ell)$ , where  $\kappa_{\ell, h_{st}}(\cdot)$  defined in (S.2) is a normalized kernel function with bandwidth parameter  $h_{st}$ .

We remark that we employ kernel smoothing twice in order to estimate the varying coefficients. In the first step, we temporally smooth the least square estimator to compute  $\tilde{\theta}_{st}^0(\iota)$ . In the second step, we further spatially smooth  $\tilde{\theta}_{st}^0(\iota)$  to compute  $\tilde{\theta}_{st}(\iota)$ . Therefore, the estimator  $\tilde{\theta}_{st}(\iota)$  has smaller variance than  $\tilde{\theta}_{st}^0(\iota)$ , since we borrow information across neighboring regions to improve the estimation efficiency. To elaborate this point, the random effect in (18) can be decomposed into three parts:  $\eta_{i,\tau,\iota}^I + \eta_{i,\tau,\iota}^{II} + \eta_{i,\tau,\iota}^{III}$ . Temporally smoothing the varying coefficient estimator removes the random fluctuations caused by  $\eta_{i,\tau,\iota}^{III}$  and the measurement error. Spatially smoothing the estimator further removes the random fluctuations caused by  $\eta_{i,\tau,\iota}^{II}$ . This in turn implies that the proposed test under the spatio-temporal design is more powerful than the one developed in Section 2 under the temporal design. Such an observation is consistent with our numerical findings in Section 5.2.

Steps 3 and 4 of Algorithm S.1 are to estimate the covariance matrix of  $(\tilde{\theta}_{st}(1), \dots, \tilde{\theta}_{st}(r))^\top$ , denoted by  $\tilde{\mathbf{V}}_{\theta, st}$ . These two steps are very similar to Steps 3 and 4 of Algorithm 1. Specifically, we first estimate the measurement errors and random effects based on the estimated varying coefficients. We next use the sandwich formula to compute the estimated covariance matrix for the initial least-square estimator. Then the estimated covariance matrix for  $\tilde{\theta}_{st}^0(\iota)$  can be derived accordingly. We use  $\tilde{\mathbf{V}}_{\theta, st}$  to denote the corresponding covariance matrix estimator.

Step 5 of Algorithm S.1 is to compute the Wald-type test statistic and its standard error estimator. Specifically, let  $\tilde{\gamma}_1(\tau, \iota)$  and  $\tilde{\gamma}_2(\tau, \iota)$  be the last two elements of  $\tilde{\theta}_{st}(\tau, \iota)$ , we have  $\widehat{\text{DE}}_{st} = \sum_{\iota=1}^r \sum_{\tau=1}^m \tilde{f}\tilde{\gamma}_1(\tau, \iota) + \tilde{\gamma}_2(\tau, \iota)g$ . We will show in Theorem 6 that  $\widehat{\text{DE}}_{st}$  is asymptotically normal. In addition, its standard error  $\widehat{\text{se}}(\widehat{\text{DE}}_{st})$  can be derived based on  $\tilde{\mathbf{V}}_{\theta, st}$ . This yields our Wald-type test statistic  $T_{st} = \widehat{\text{DE}}_{st} / \widehat{\text{se}}(\widehat{\text{DE}}_{st})$ . We reject the null hypothesis if  $T_{st}$  exceeds the upper  $\alpha$ th quantile of a standard normal distribution.

We next describe our estimation and testing procedures for IE. The method is very similar to the one discussed in Section 2.4. We sketch an outline of the algorithm to save space. Details are presented in S.2 of Section S.1 of the supplementary document. Specifically, we first plug in the set of smoothed estimators  $\tilde{f}\tilde{\Theta}_{st}(\tau, \iota)g_{\tau,\iota}$  and  $\tilde{f}\tilde{\theta}_{st}(\tau, \iota)g_{\tau,\iota}$  for  $\tilde{f}\Theta(\tau, \iota)g_{\tau,\iota}$  and  $\tilde{f}\theta(\tau, \iota)g_{\tau,\iota}$  to compute  $\widehat{\text{IE}}_{st}$ , the plug-in estimator of  $\text{IE}_{st}$ . We next estimate the measurement errors and random effects and then apply the parametric bootstrap method to compute the bootstrap statistics  $\widehat{\text{IE}}_{st}^b$ . Finally, we reject  $H_0^{IE}$  if  $\widehat{\text{IE}}_{st}$  exceeds the upper  $\alpha$ th empirical quantile of  $\widehat{\text{IE}}_{st}^b$ .

To conclude this section, we remark that in Sections 2 and 3, we focus on testing one-sided hypotheses for the direct and indirect effects. However, the proposed method can be easily extended to test two-sided hypotheses as well.

#### 4. Theoretical Analysis

In this section, we systematically investigate the asymptotic properties of the proposed estimators and test statistics in L-TVCDP and derive the convergence rates of our causal estimands in NN-TVCDP. We also explore the benefits of employing the switchback design and study the theoretical properties of our estimator in the spatio-temporal dependent experiments.

Firstly, we impose the following regularity assumptions for the temporal dependent experiments using L-TVCDP.

**ASSUMPTION 3.** *The kernel function  $K(\cdot)$  is a symmetric probability density function on  $[-1, 1]$  and is Lipschitz continuous.*

**ASSUMPTION 4.** *The covariate  $\mathbf{Z}_i$ s are i.i.d.; for  $1 \leq \tau \leq m$ ,  $\mathbb{E}(\mathbf{Z}_{i,\tau}^\top \mathbf{Z}_{i,\tau}) \in \mathbb{M}^{p \times p}$  is invertible; all components of  $\theta(t)$  have bounded and continuous second derivatives with respect to  $t$ .*

**ASSUMPTION 5.** *There exists  $0 < q < 1$  such that the absolute values of eigenvalues of  $\Phi(\tau)$  are smaller than  $q$ , and there exist some constants  $M_\Gamma$  and  $M_\beta$  such that  $k\Gamma(\tau)k_\infty \leq M_\Gamma$  and  $k\beta(\tau)k_\infty \leq M_\beta$ .  $f\beta(\tau)g_{2 \leq \tau \leq m}$ ,  $f\Phi(l)g_{2 \leq l \leq m-1}$ , and  $f\Gamma(k)g_{1 \leq k \leq m-1}$  must not be all zero.  $\Theta(\tau)$  has a continuous second-order partial derivative.*

Assumption 3 is mild as the kernel  $K(\cdot)$  is user-specified. Assumption 4 has been commonly used in the literature on varying coefficient models (see e.g., Zhu et al., 2014). Assumption 5 ensures that the time series is stationary, since  $\Phi(\tau)$  is the autoregressive coefficient. It is commonly imposed in the literature on time series analysis (Shumway and Stoffer, 2010).

Before presenting the theoretical properties of the proposed method for L-TVCDP, we introduce some notation. For  $1 \leq \tau_1, \tau_2 \leq m$ , define  $\Sigma_y$  and  $\Sigma_\eta$  to be the  $m \times m$  matrices  $f\Sigma_y(\tau_1, \tau_2)g_{\tau_1, \tau_2}$  and  $f\Sigma_\eta(\tau_1, \tau_2)g_{\tau_1, \tau_2}$ , respectively. We define

$$\mathbf{V}_{\hat{\theta}} = (\mathbb{E}\mathbf{Z}_i^\top \mathbf{Z}_i)^{-1} \mathbb{E}(\mathbf{Z}_i^\top \Sigma_y \mathbf{Z}_i) (\mathbb{E}\mathbf{Z}_i^\top \mathbf{Z}_i)^{-1} \quad \text{and} \quad \mathbf{V}_{\tilde{\theta}} = (\mathbb{E}\mathbf{Z}_i^\top \mathbf{Z}_i)^{-1} \mathbb{E}(\mathbf{Z}_i^\top \Sigma_\eta \mathbf{Z}_i) (\mathbb{E}\mathbf{Z}_i^\top \mathbf{Z}_i)^{-1}$$

as the asymptotic covariance matrices of  $\hat{\theta}$  and  $\tilde{\theta}$ , respectively. Let  $\mathbf{V}_{\hat{\theta}}(\tau, \tau)$  and  $\mathbf{V}_{\tilde{\theta}}(\tau, \tau)$  denote the submatrices of  $\mathbf{V}_{\hat{\theta}}$  and  $\mathbf{V}_{\tilde{\theta}}$  that correspond to the asymptotic covariance matrix of  $\hat{\theta}$  and  $\tilde{\theta}$ , respectively. We first compare the mean squared error (MSE) of the OLS estimator  $\hat{\theta}(\tau)$  against that of the smoothed estimator  $\tilde{\theta}(\tau)$  based on L-TVCDP.

**PROPOSITION 2.** *Suppose  $\lambda_{\min}(\mathbf{V}_{\hat{\theta}}(\tau, \tau))$  and  $\lambda_{\min}(\mathbf{V}_{\tilde{\theta}}(\tau, \tau))$  are uniformly bounded away from zero for any  $\tau$ . Under Assumptions 3 and 4, we have*

$$\sum_{\tau=1}^m \text{MSE}(\hat{\theta}(\tau)) \leq n^{-1} \text{trace}(\mathbf{V}_{\hat{\theta}}), \quad \sum_{\tau=1}^m \text{MSE}(\tilde{\theta}(\tau)) \leq n^{-1} \text{trace}(\mathbf{V}_{\tilde{\theta}}) + O(mh^4 + m^{-1}).$$

Proposition 2 has an important implication. Both  $\text{trace}(\mathbf{V}_{\hat{\theta}})$  and  $\text{trace}(\mathbf{V}_{\tilde{\theta}})$  are of the order of magnitude  $O(m)$ . When  $m \asymp \frac{\rho}{n}$  or  $h^4 \asymp n^{-1}$ , the squared bias of  $\tilde{\theta}$  may dominate its variance. Hence, the OLS estimator  $\hat{\theta}$  may achieve a smaller MSE. When  $m \asymp \frac{\rho}{n}$  and  $h^4 = O(n^{-1}m)$ , the two MSEs are of the same order of magnitude and it remains unclear which one is smaller. When  $m \asymp \frac{\rho}{n}$  and  $h^4 = o(n^{-1})$ , the variance of  $\tilde{\theta}$  dominates its squared bias. Moreover,  $\Sigma_y - \Sigma_\eta$  is strictly positive definite, so is  $\mathbf{V}_{\hat{\theta}} - \mathbf{V}_{\tilde{\theta}}$ . As a result,  $\tilde{\theta}$  achieves a smaller MSE. In our applications,  $m$  is moderately large and the condition  $m \asymp \frac{\rho}{n}$  is likely to be satisfied. With properly chosen bandwidth, we expected the smoothed estimator achieves a smaller MSE.

Secondly, we present the limiting distributions of  $\hat{\theta}(\tau)$  and  $\tilde{\theta}(\tau)$  and prove the validity of our test for DE based on L-TVCDP.

**THEOREM 1.** *Suppose  $\lambda_{\min}(\mathbf{V}_{\hat{\theta}}(\tau, \tau))$  and  $\lambda_{\min}(\mathbf{V}_{\tilde{\theta}}(\tau, \tau))$  are uniformly bounded away from zero for any  $\tau$ . Under Assumptions 1, 3 and 4, for any  $(d+2)$ -dimensional vectors  $\mathbf{a}_{n,1}, \mathbf{a}_{n,2}$ , with unit  $\ell_2$  norm,*

$$(i) \quad \frac{\rho}{n} \mathbf{a}_{n,1}^\top f\hat{\theta}(\tau) - \theta(\tau)g / \sqrt{\mathbf{a}_{n,1}^\top \mathbf{V}_{\hat{\theta}}(\tau, \tau) \mathbf{a}_{n,1}} \xrightarrow{d} N(0, 1) \quad \text{as } n \rightarrow \infty \text{ for any } \tau;$$

$$(ii) \quad \text{Suppose } m \rightarrow \infty, h \rightarrow 0, \text{ and } hm \rightarrow 1 \text{ as } n \rightarrow \infty. \text{ Then } \frac{\rho}{n} \mathbf{a}_{n,2}^\top f\tilde{\theta}(\tau) - \theta(\tau)g / \sqrt{\mathbf{a}_{n,2}^\top \mathbf{V}_{\tilde{\theta}}(\tau, \tau) \mathbf{a}_{n,2}} \xrightarrow{d} N(b_n, 1) \text{ as } n \rightarrow \infty \text{ for any } \tau, \text{ where the bias } b_n = O\left(\frac{\rho}{n} h^2 + \frac{\rho}{n} m^{-1}\right).$$

(iii) Suppose  $h = o(n^{-1/4})$ ,  $m \stackrel{P}{\rightarrow} \infty$  and the sum of all elements in  $m^{-2}\mathbf{V}_{\tilde{\gamma}}$  is bounded away from zero where  $\mathbf{V}_{\tilde{\gamma}}$  denotes the submatrix of  $\mathbf{V}_{\tilde{\theta}}$  which corresponds to the asymptotic covariance matrix of  $\tilde{\theta}$ . Then for the hypotheses (2), under  $H_0^{DE}$ ,  $\mathbb{P}(\widehat{DE}/\widehat{se}(\widehat{DE}) > z_\alpha) = \alpha + o(1)$ ; under  $H_1^{DE}$ ,  $\mathbb{P}(\widehat{DE}/\widehat{se}(\widehat{DE}) > z_\alpha) \rightarrow 1$ , where  $z_\alpha$  denotes the upper  $\alpha$ th quantile of a standard normal distribution.

Theorem 1 has several important implications. First, the bias of the smoothed estimator  $\tilde{\theta}$  decays with  $m$ . In cases where  $m$  is fixed, the kernel smoothing step is not preferred as it will result in an asymptotically biased estimator. Second, each  $\tilde{\theta}(\tau)$  converges at a rate of  $O_p(n^{-1/2})$  under the assumption that  $\lambda_{\min}(\mathbf{V}_{\tilde{\theta}}(\tau, \tau))$  is bounded away from zero. The rate  $O_p(n^{-1/2}m^{-1/2})$  cannot be achieved despite that we have a total of  $nm$  observations, since the random errors  $f_{e_\tau}g_\tau$  are not independent. We also remark that in the extreme case where  $f_{e_\tau}g_\tau$  are independent, we can set  $h \propto (nm)^{-1/5}$  and  $\tilde{\theta}(\tau)$  attains the classical nonparametric convergence rate  $O_p((nm)^{-2/5})$ . Third, since  $\mathbf{V}_{\tilde{\theta}} = \mathbf{V}_{\tilde{\theta}}$  is strictly positive, this similarly implies that the smoothed estimator is more efficient when  $b_n = o(1)$ , or equivalently,  $h = o(n^{-1/4})$  and  $m \stackrel{P}{\rightarrow} \infty$ . Finally, in the proof of Theorem 1, we show that the covariance estimator  $\widehat{\mathbf{V}}_{\tilde{\theta}}$  is consistent. This together with asymptotic distribution of  $\tilde{\theta}$  yields the consistency of our test in (iii).

Thirdly, we present the validity of the proposed parametric bootstrap procedure for IE under the temporal alternation design based on L-TVCDP.

**THEOREM 2.** *Suppose that there is some constant  $0 < c_1 < 1$  such that  $c_1 \leq \mathbb{E}k_{\varepsilon_\tau, S}^2$  and  $\mathbb{E}e_\tau^2 \leq c_1^{-1}$  for all  $1 \leq \tau \leq m$ . Suppose that  $h = o(n^{-1/4})$ ,  $m \leq n^{c_2}$  for some  $1/2 \leq c_2 < 3/2$  and  $mh \rightarrow 1$ . Then under the assumptions in Theorem 1 and Assumption 5, with probability approaching 1, we have*

$$\sup_z \mathbb{P}(\widehat{IE} - IE \leq z) - \mathbb{P}(\widehat{IE}^b - \widehat{IE} \leq z | \text{Data}) \leq C(n^{-1/2}h^2 + n^{-1/8}),$$

where  $C$  is some positive constant.

We have several remarks. The derivation of Theorem 2 is non-trivial when  $m$  diverges with  $n$ . Specifically, since  $\widehat{IE}$  is a very complicated function of the estimated varying coefficients (see Equation (14)), its limiting distribution is not well-defined. To prove Theorem 2, we derive a nonasymptotic error bound on the difference between the distribution of  $\widehat{IE}$  and that of the bootstrap statistics conditional on the data. As a result, it ensures that the type-I error can be well-controlled and the power approaches one. Please refer to the proof of Theorem 2 in the supplementary document for details. Finally, we require  $m$  to diverge with  $n$  at certain rate. In settings with a small or fixed  $m$ , one can apply the proposed bootstrap procedure to the unsmoothed estimator  $\hat{\theta}$ . The resulting test procedure remains valid regardless of whether  $m$  is fixed or not.

Fourthly, we illustrate the advantage of employing the switchback design in the presence of temporal random effects. As commented in the introduction, the switchback design assigns different treatments at adjacent time points  $A_{i,1} = 1, A_{i,2} = A_{i,3} = \dots = A_{i,2t-1} = 1, A_{i,2t}$ , whereas the alternating-day design assigns fixed treatment  $A_{i,1} = A_{i,2} = A_{i,3} = \dots = A_{i,2t-1} = A_{i,2t}$  within each day for any  $i$  and  $t$ . In the switchback design, the random effects at adjacent time points can cancel with each other when estimating the causal effect, yielding a more efficient estimator. To elaborate this point, we compare the mean square errors of the proposed estimators under the switchback design against those under an alternating-day design where the new and old policies are daily switched back and forth. To simplify the analysis, we focus on the case where the state is one-dimensional and assume the treatment effect estimators are constructed based on the unsmoothed OLS estimators (see Section S.12.3 for details). Let  $\text{MSE}(\widehat{DE}_{sb})$  and  $\text{MSE}(\widehat{DE}_{ad})$  denote the mean squared errors of DE estimators under the switchback design and the alternating-day design, respectively.

**THEOREM 3.** *Suppose that the state is one-dimensional,  $\Sigma_\eta(\tau_1, \tau_2)$  is nonnegative for any  $\tau_1$  and  $\tau_2$  and Assumptions 1 and 4 hold. When  $f\Phi(\tau)g_\tau$  and  $f\Gamma(\tau)g_\tau$  are of the same signs, respectively, i.e. for any  $\tau_1, \tau_2$ ,  $\Phi(\tau_1)\Phi(\tau_2) \geq 0$  and  $\Gamma(\tau_1)\Gamma(\tau_2) \geq 0$ , then as  $n \rightarrow \infty$ , we have*

$$n\text{MSE}(\widehat{DE}_{sb}) = n\text{MSE}(\widehat{DE}_{ad}) + o(1),$$

where the equality holds only when  $\Sigma_\eta(j, k) = 0$  for any  $j, k$  such that  $jj - kk = 1, 3, 5, \dots$

To ensure that DE achieves a much smaller MSE under the switchback design, we only require that the random effects are non-negatively correlated and that the correlation  $\Sigma(j, k)$  is nonzero for some  $j = k = 1, 3, 5, \dots$ . These conditions are automatically satisfied when the random effects are positively correlated. We next provide a close-formed expression for the ratio of the two MSEs under an AR(1) noise structure and the constraint that  $\Gamma(1) = \Gamma(2) = \dots = \Gamma(m-1) = 0$ .

**COROLLARY 1.** *Suppose that for any  $1 \leq \tau_1, \tau_2 \leq m$ ,  $\Sigma_e(\tau_1, \tau_2) = c\rho^{|\tau_1 - \tau_2|}$  for some constant  $c > 0$ . Then under assumptions of Theorem 3, when  $\Gamma(1) = \Gamma(2) = \dots = \Gamma(m-1) = 0$ , we have as  $n, m \rightarrow \infty$ ,*

$$\frac{\text{MSE}(\widehat{DE}_{sb})}{\text{MSE}(\widehat{DE}_{ad})} = \frac{(1 - \rho)^2}{(1 + \rho)^2} + o(1).$$

It can be seen from Corollary 1 that the larger the  $\rho$ , the smaller the variance ratio. In particular, when  $\rho = 0.5$ , MSE of DE under the switchback design is approximately 9 times smaller than that under the alternating-day design. We next consider IE.

**THEOREM 4.** *Suppose  $m = 2$ . Under Assumptions 1 and 4, we have*

$$n \text{fMSE}(\widehat{IE}_{ad}) - \text{MSE}(\widehat{IE}_{sb}) = o(1).$$

Theorem 4 suggests that the IE estimators under the two designs have comparable MSEs. This together with Theorem 3 underscores the superiority of the switchback design, particularly when  $m = 2$ . However, as  $m$  exceeds 2, determining the closed-form expression for  $\text{MSE}(\widehat{IE})$  becomes exceedingly complex, making it challenging to directly compare the two designs. Addressing this complexity and extending the comparison for cases where  $m > 2$  is a task we reserve for future research.

Fifth, we establish the convergence rates of the estimated DE and IE for NN-VCDP.

**THEOREM 5.** *Suppose that  $f_{\varepsilon_{\tau S}}$  is Lipschitz, meaning that for any  $\tau$ , there exists a constant  $L_f > 0$  such that  $|f_{\varepsilon_{\tau S}}(x) - f_{\varepsilon_{\tau S}}(y)| \leq L_f \|x - y\|_{k_2}$ , where  $\|\cdot\|_{k_2}$  represents the Frobenius norm. Additionally, assume that the NN-based learners satisfy  $\mathbb{E}f\widehat{G}_a(\tau, S_\tau) - G_a(\tau, S_\tau) \leq \Delta_1^2(n, m)$  and  $\mathbb{E}f\widehat{g}_a(\tau, S_\tau^{a_1}) - g_a(\tau, S_\tau^{a_1}) \leq \Delta_2^2(n, m)$ , where  $a \geq 2$ ,  $l, g$  and  $\Delta_1(n, m)$  and  $\Delta_2(n, m)$  are specific functions. The density estimator should fulfill  $\|\int_x f_{\varepsilon_{\tau S}}(x) - \widehat{f}_{\varepsilon_{\tau S}}(x) dx\| = O_p(\Delta_3(n, m))$  for some function  $\Delta_3$ . Both  $g_a$  and  $\widehat{g}_a$  must be uniformly bounded. Moreover, the ratio of the density function of the potential state  $S_\tau^a$  to the density of the observed state  $S_\tau$  must be bounded by  $\frac{1}{\rho}$  for any  $\tau$  and  $a$ . Then, as  $\min(n, m) \rightarrow \infty$ , we obtain the following convergence results:*

$$\begin{aligned} \widehat{DE} - DE &= O_p \left( m \frac{\rho}{\omega} \Delta_2(n, m) + m^2 \Delta_1(n, m) + m^2 L_f \frac{\rho}{\omega} \Delta_3(n, m) + \frac{m}{n} \sqrt{\log(nm)} \right), \\ \widehat{IE} - IE &= O_p \left( m \frac{\rho}{\omega} \Delta_2(n, m) + m^2 \Delta_1(n, m) + m^2 L_f \frac{\rho}{\omega} \Delta_3(n, m) + \frac{m}{n} \sqrt{\log(nm)} \right). \end{aligned}$$

Since the convergence rates of NN-based learners have been widely studied in the literature (see e.g., Shen et al., 2019; Schmidt-Hieber, 2020; Shen et al., 2022; Yan and Yao, 2023), these results can be used to establish the convergence rates of  $\widehat{G}_a$  and  $\widehat{g}_a$ .

Finally, we impose the following regularity assumptions for the proposed tests in spatio-temporal dependent experiments based on L-STVCDP.

**ASSUMPTION 6.** *For any  $\tau, \iota$ ,  $\mathbb{E}(Z_{i, \tau, \iota}^\top Z_{i, \tau, \iota})$  is invertible;  $\theta(\tau, \iota)$ ,  $\Sigma_{\eta^t}(\tau_1, \tau_2, \iota_1, \iota_2)$ ,  $\Sigma_{\eta^{\iota_1}}(\tau_1, \iota_1, \tau_2)$ , and  $\Sigma_{\eta^{\iota_1 \iota_2}}(\tau_1, \iota_1, \iota_2)$  have bounded and continuous second-order derivatives.*

**ASSUMPTION 7.** *There exists  $q < 1$  such that the absolute values of eigenvalues of  $\Phi(\tau, \iota)$  are smaller than  $q$ . In addition, there exist  $M_\Gamma$  and  $M_\beta < 1$  such that  $k\Gamma_1(\tau, \iota) + \Gamma_2(\tau, \iota) \leq M_\Gamma$  and  $k\beta(\tau, \iota) \leq M_\beta$ .  $\Theta(\tau, \iota)$  has a bounded and continuous second-order derivative.*

With these assumptions, we present the asymptotic properties of our DE and IE estimators and their associated test statistics for the spatio-temporal dependent experiments based on L-STVCDP. Define

$$\mathbf{V}_{\widehat{\theta}_{st}}(\tau_1, \iota_1, \tau_2, \iota_2) = \text{f}\mathbb{E}Z_{i, \tau_1, \iota_1} Z_{i, \tau_1, \iota_1}^\top \mathbf{g}^{-1} \text{f}\mathbb{E}fZ_{i, \tau_2, \iota_2} Z_{i, \tau_1, \iota_1}^\top \Sigma_{\eta^t}(\tau_1, \iota_1, \tau_2, \iota_2) \text{g}\text{f}\mathbb{E}Z_{i, \tau_2, \iota_2} Z_{i, \tau_2, \iota_2}^\top \mathbf{g}^{-1}$$

as the asymptotic covariance between  $\frac{\rho}{n} \widehat{\theta}_{st}(\tau_1, \iota_1)$  and  $\frac{\rho}{n} \widehat{\theta}_{st}(\tau_2, \iota_2)$ .

THEOREM 6. Suppose  $\lambda_{\min}(\mathbf{V}_{\tilde{\theta}_{st}})$  is bounded away from zero. Under Assumptions 2, 3 and 6, for any set of  $(d+2)$ -dimensional vectors  $f_{B_{\tau,\ell}}g_{\tau,\ell}$ , we have as  $n, m, r \rightarrow \infty$ ,  $h, h_{st} \rightarrow 0$  and  $mh, rh_{st} \rightarrow 1$  that

(i) For any set of  $(d+2)$ -dimensional vectors  $f_{B_{\tau,\ell}}g_{\tau,\ell}$  with  $\sum_{\tau_1, \tau_2, \ell_1, \ell_2} B_{\tau_1, \ell_1}^\top \mathbf{V}_{\tilde{\theta}_{st}}(\tau_1, \ell_1, \tau_2, \ell_2) B_{\tau_2, \ell_2} = c \sum_{\tau, \ell} k_{B_{\tau, \ell}}^2$  for some constant  $c > 0$ , we have

$$\frac{\rho_{\bar{n}} \sum_{\tau, \ell} [B_{\tau, \ell}^\top f_{\tilde{\theta}_{st}}(\tau, \ell) - \theta_{st}(\tau, \ell)g] / \sqrt{\sum_{\tau_1, \tau_2, \ell_1, \ell_2} B_{\tau_1, \ell_1}^\top \mathbf{V}_{\tilde{\theta}_{st}}(\tau_1, \ell_1, \tau_2, \ell_2) B_{\tau_2, \ell_2}}}{\rho_{\bar{n}}} \xrightarrow{d} N(b_{n, st}, 1),$$

where the bias  $b_{n, st} = O(\rho_{\bar{n}} h^2 + \rho_{\bar{n}} h_{st}^2 + \rho_{\bar{n}} m^{-1} + \rho_{\bar{n}} r^{-1})$ .

(ii) Suppose  $h, h_{st} = o(n^{-1/4})$  and  $m, r \rightarrow \infty$ . Then for the hypotheses (16),  $\mathbb{P}(\widehat{DE}_{st}/\widehat{se}(\widehat{DE}_{st}) > z_\alpha) = \alpha + o(1)$  under  $H_0^{DE}$  and  $\mathbb{P}(\widehat{DE}_{st}/\widehat{se}(\widehat{DE}_{st}) > z_\alpha) \rightarrow 1$  under  $H_1^{DE}$ .

THEOREM 7. Suppose that there are some constants  $0 < c_1 < 1$  such that  $c_1 \leq \mathbb{E}\varepsilon_{\tau, \ell, S}^2, \mathbb{E}e_{\tau, \ell}^2 \leq c_1^{-1}$  for all  $1 \leq \tau \leq m, 1 \leq \ell \leq r$ , and that  $h, h_{st} = o(n^{-1/4})$ ,  $m, r \rightarrow \infty$  and  $mr \rightarrow n^{c_2}$  for some constant  $c_2 < 3/2$ . Then under Assumptions of Theorem 6 and Assumption 7, with probability approaching 1,

$$\sup_z \mathbb{P}(\widehat{IE}_{st} - IE_{st} \leq z) - \mathbb{P}(\widehat{IE}_{st}^b - \widehat{IE}_{st} \leq z | \text{Data}) \leq C(\rho_{\bar{n}} h^2 + \rho_{\bar{n}} h_{st}^2 + \rho_{\bar{n}} m^{-1} + \rho_{\bar{n}} r^{-1} + n^{-1/8}), \quad (20)$$

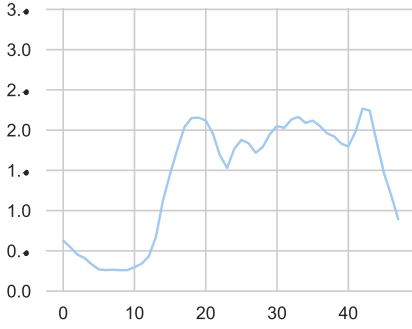
where  $C$  is some positive constant.

Theorem 6 establishes the limiting distribution of the proposed DE estimator for the spatio-temporal dependent experiments. Similar to Proposition 2, we can show that the smoothed estimator is more efficient when  $m, r \rightarrow \infty$  and  $h^4, h_{st}^4 = o(n^{-1})$ . In addition, Theorem 7 allows both  $m$  and  $r$  to be either fixed, or diverge with  $n$ , and is thus applicable to a wide range of applications.

## 5. Real data based simulations

### 5.1. Temporal alternation design

In this section, we conduct Monte Carlo simulations to examine the finite sample properties of the proposed test statistics based on L-TVCDP and L-STVCDP models. To generate data under the temporal alternation design, we design two simulation environments based on two real datasets obtained from Didi Chuxing. The first dataset is collected from a given city A from Dec. 5th, 2018 to Jan. 13th, 2019. Thirty-minutes is defined as one time unit. The second dataset is from another city B, from May 17th, 2019 to June 25th, 2019. One-hour is defined as one time unit. Both contain data for 40 days. Due to privacy, we only present scaled metrics in this paper. Figure 1 depicts the trend of some business metrics over time across 40 different days. These metrics include drivers' total income, the number of requests and drivers' total online time. Among them, the first quantity is our outcome of interest and the last two are considered as the state variables to characterize the demand and supply networks. As expected, these quantities show a similar pattern, achieving the largest values at peak time.

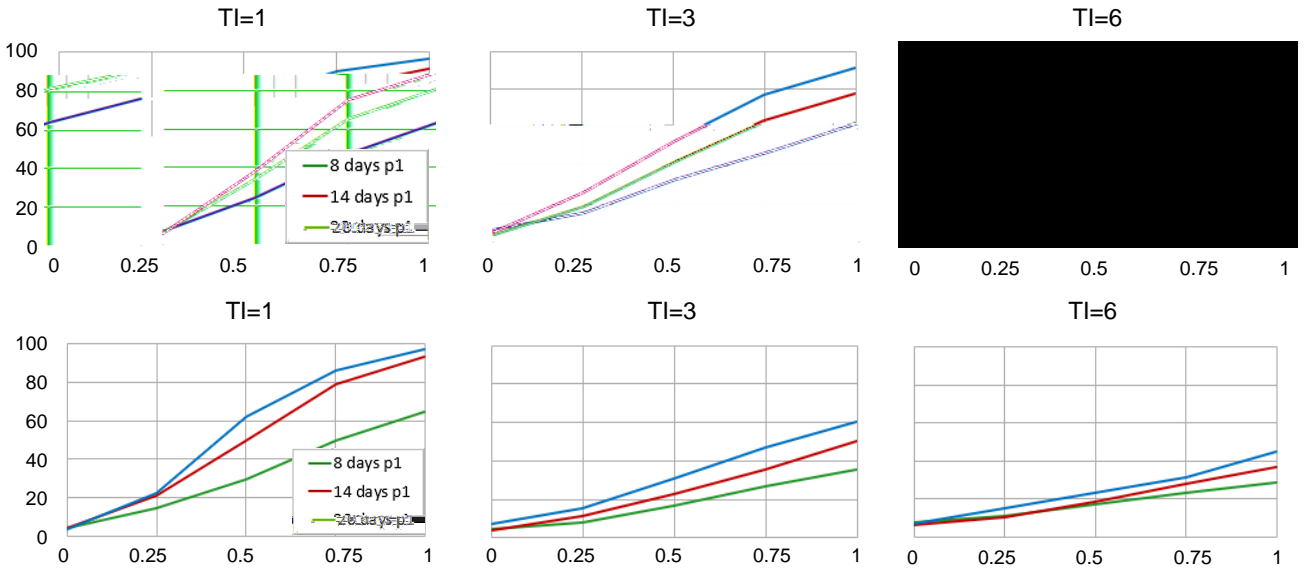


**Fig. 1.** Scaled business metrics from City A (the first row) and City B (the second row) across 40 days, including drivers' total income, the numbers of requests and drivers' total online time.

We next discuss how to generate synthetic data based on the real datasets. The main idea is to fit the proposed L-TVCDP models to the real dataset and apply the parametric bootstrap to simulate the data. Let  $\tilde{\beta}_0(\tau)$ ,  $\tilde{\beta}(\tau)$ ,  $\tilde{\phi}_0(\tau)$ , and  $\tilde{\Phi}(\tau)$  denote the smoothed estimators for  $\beta_0(\tau)$ ,  $\beta(\tau)$ ,  $\phi_0(\tau)$  and  $\Phi(\tau)$ , respectively. We set  $\tilde{\gamma}(\tau)$  and  $\tilde{\Gamma}(\tau)$  to  $(\delta/100) \left( \sum_{i,\tau} Y_{i,\tau}/nm \right)$  and  $(\delta/100) \left( \sum_{i,\tau} S_{i,\tau}/nm \right)$ , respectively. As such, the parameter  $\delta$  controls the degree of the treatment effects. Specifically, the null holds if  $\delta = 0$  and the alternative holds if  $\delta > 0$ . It corresponds to the increase relative to the outcome (state). We next generate the policies according to the temporal alternation design and simulate the responses and states based on the fitted model. Let TI denote the time span we implement each policy under the alternation design. For instance, if  $\text{TI} = 3$ , then we first implements one policy for three hours, then switch to the other for another three hours and then switch back and forth between the two policies. We consider three choices of  $n \in \{8, 14, 20\}$ , five choices of  $\delta \in \{0, 0.25, 0.5, 0.75, 1\}$  and three choices of  $\text{TI} \in \{1, 3, 6\}$ . This corresponds to a total of 45 cases. The bandwidth is set  $h = Cn^{-1/3}$ , where  $C$  is selected by the 5-fold cross validation method.

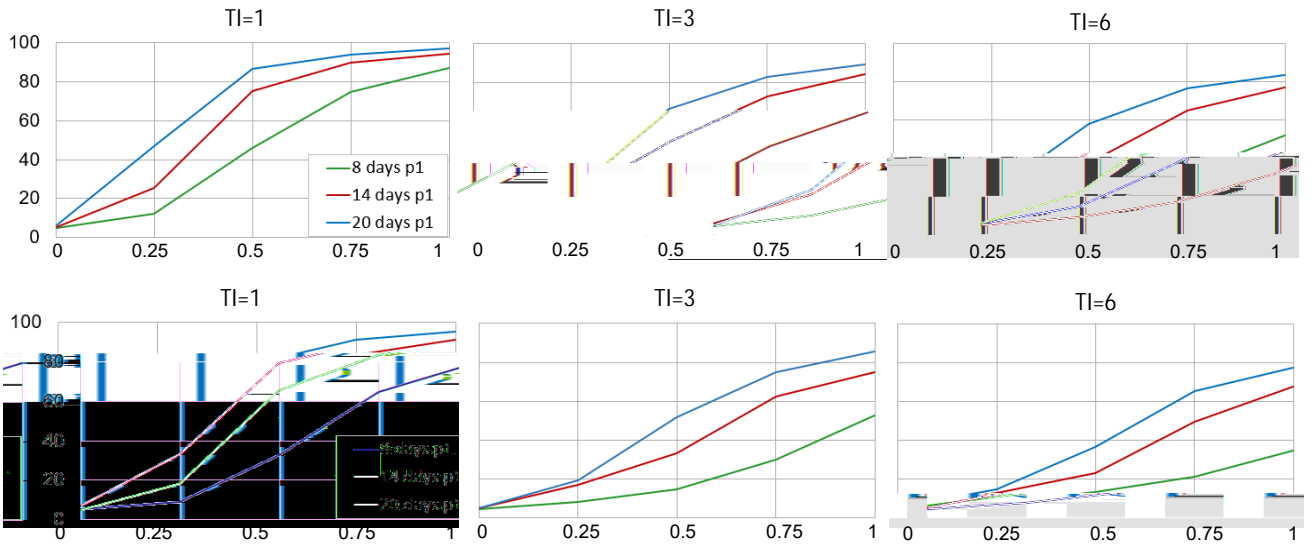
In Figure 2, we depict the empirical rejection probabilities of the proposed test for DE, aggregated over 400 simulations, for all combinations. It can be seen that our test controls the type-I error and its power increases as  $\delta$  increases. In addition, the empirical rejection rates decreases as TI increases. This phenomenon suggests that the more frequently we switch back and forth between the two policies, the more powerful the resulting test. It is due to the positive correlation between adjacent observations. To elaborate, consider the extreme case where we switch policies at each time. The policies assigned at any two adjacent time points are different. As such, the random effect cancels with each other, yielding an efficient estimator. We conduct some additional simulations using the numbers of answered requests and finished requests of cities A and B as responses (see Figure S.2 in the supplement). Results are very similar and are reported in Figures S.3–S.4 in the supplementary document. See also Tables S.1–S.2 in the supplementary document.





**Fig. 2.** Simulation results for L-TVCDP: empirical rejection rates (expressed as percentages) of the proposed test for DE under different combinations of  $(n, \delta, TI)$  and types of outcomes. Synthetic data are simulated based on the real dataset from city A (the first row) and city B (the second row).

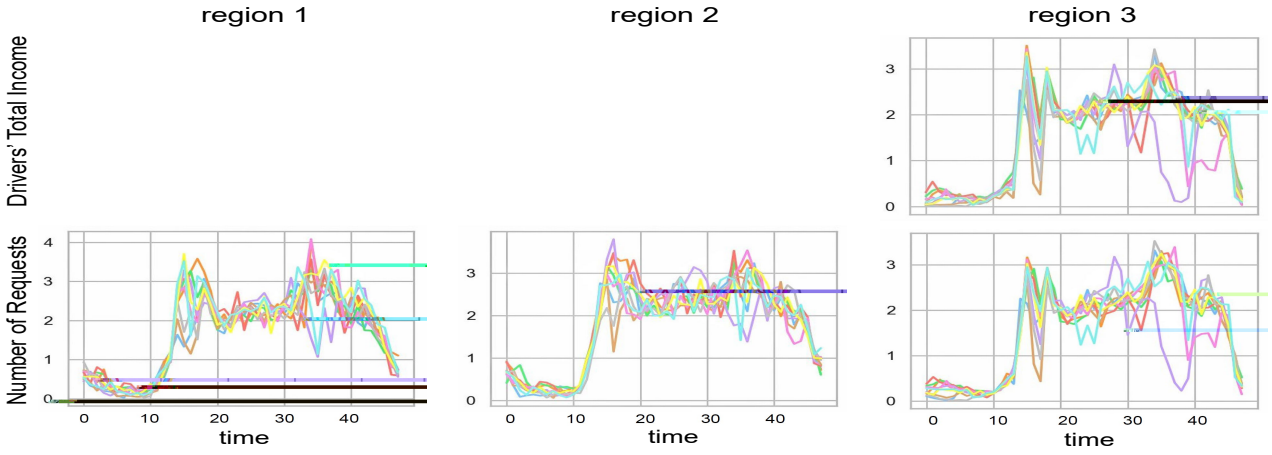
To infer IE, we set the outcome to drivers’ total online income. The empirical rejection probabilities of the proposed test for IE are reported in Figure 3. Results are aggregated over 400 simulations. Similarly, the proposed test is consistent. Its power increases with the sample size and  $\delta$ . In addition, its power under  $TI = 1$  is much larger than those under  $TI = 3$  or 6. This suggests that we shall switch back and forth between the two policies as frequently as possible to maximize the power property of the test (see also Tables S.3–S.4 in Supplementary document).



**Fig. 3.** Simulation results for L-TVCDP: empirical rejection rates (expressed as percentages) of the proposed test for IE under different combinations of  $(n, \delta, TI)$ . Synthetic data are simulated based on the real dataset from city A (the first row) and city B (the second row).

### 5.2. Spatio-temporal alternation design

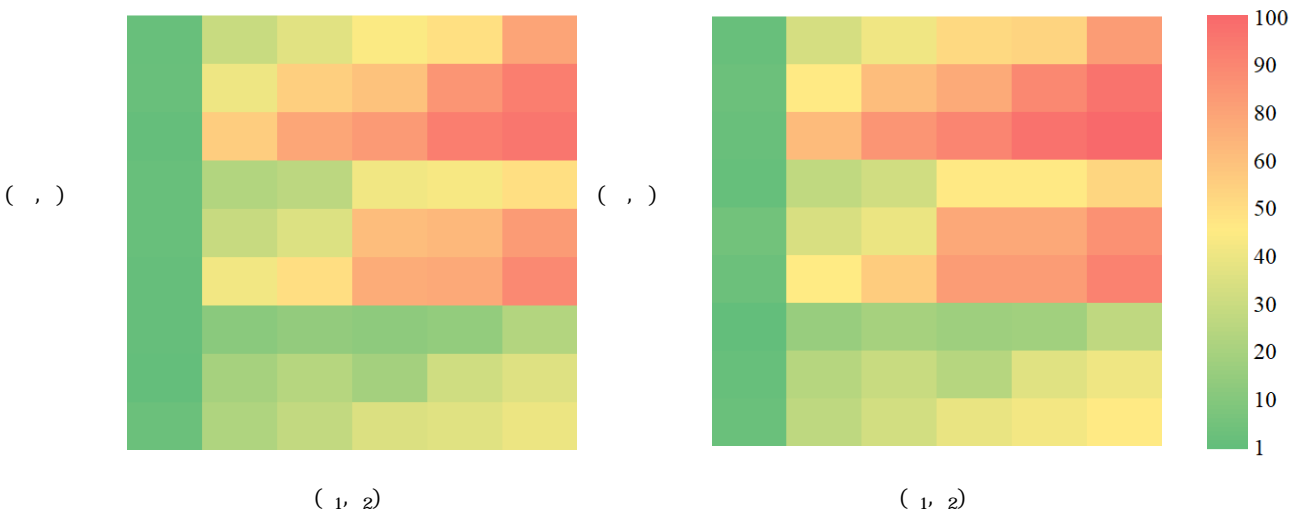
To generate data under the spatio-temporal alternation design, we create a simulation environment based on the real dataset from city A. We divide the city into 10 non-overlapping regions. We plot these variables associated with 3 particular regions, over the first 10 days in Figure 4. It can be seen that although the daily trends differ across regions, the state and the response are highly correlated.



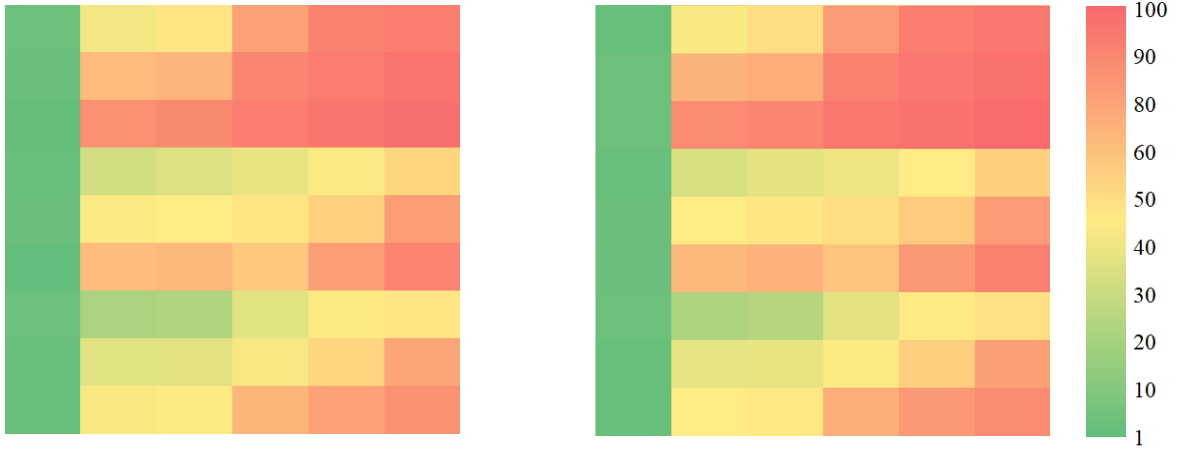
**Fig. 4.** Number of call requests and drivers' total income across different regions and days. The values are scaled for privacy concerns.

We fit the proposed models in (18) to the real dataset to estimate the varying coefficients and the variances of the random errors. Then we manually set the treatment effects  $\hat{\gamma}(\tau, \iota)$  and  $\hat{\Gamma}(\tau, \iota)$  to  $(\delta_1/100) \left( \sum_{i=1}^n \sum_{\tau=1}^m Y_{i,\tau,\iota} / nm \right)$  and  $(\delta_2/100) \left( \sum_{i=1}^n \sum_{\tau=1}^m S_{i,\tau,\iota} / nm \right)$  for some constants  $\delta_1$  and  $\delta_2 > 0$ . We consider both the temporal and spatio-temporal alternation designs, and simulate the data via parametric bootstrap.

We also consider three choices of  $n \in \{8, 14, 20\}$ , three choices of TI  $\in \{1, 3, 6\}$  and three choices of  $\delta_1, \delta_2 \in \{0, 0.5, 1\}$ . This yields a total of 81 combinations under each design. The rejection probabilities of the proposed tests for DE and IE tests are reported in Figures 5 and 6 (see also Tables S.5 and S.6 in the supplementary document). It can be seen that the type I error rates of the proposed test are close to the nominal level under both designs. More importantly, the power under spatio-temporal alternation design is higher than that of temporal alternation design in all cases. The reason is twofold. First, under the spatio-temporal design, we independently randomize the initial policy for each region, and adjacent regions may receive different policies. Observations across adjacent areas are likely to be positively correlated. As such, the variance of the estimated treatment effects will be smaller than that under the temporal design where all regions receive the same policy at each time. Second, we employ kernel smoothing twice when computing  $\widehat{DE}_{st}$  and  $\widehat{IE}_{st}$ , as discussed in Section 3. This results in a more efficient estimator. In addition, compared with the results in Tables S.1 and S.3, it can be seen that the test that focuses on the entire city has better power property than the one that considers a particular region in general. Finally, the power decreases with TI and increases with  $n, \delta_1$  and  $\delta_2$ .



**Fig. 5.** Simulation results for L-STVCDP: the empirical rejection probabilities of the proposed test test for DE



**Fig. 6.** Simulation results for L-STVCDP: the empirical rejection probabilities of the proposed test for IE under the temporal alternation design (left panel) and the spatio-temporal alternation design (right panel).

## 6. Real data analysis

In this section, we apply the proposed tests based on L-TVCDP and L-STVCDP to a number of real datasets from Didi Chuxing to examine the treatment effects of some newly developed order dispatch and vehicle reposition policies. Due to privacy, we do not publicize the names of these policies.

We first consider four data sets collected from four online experiments under the temporal alternation design. All the experiments last for 14 days. Policies are executed based on alternating half-hourly time intervals. We denote the cities, in which these experiments take place, as  $C_1, C_2, C_3$ , and  $C_4$  and their corresponding policies as  $S_1, S_2, S_3$ , and  $S_4$ , respectively. For each policy, we are interested in its effect on three key business metrics, including drivers' total income, the answer rate, and the completion rate. Similar to Section 5.1, we use the number of call orders and drivers' total online time to construct the time-varying state variables.

All the new policies are compared with some baseline policies in order to evaluate whether they improve some business outcomes. Specifically, in city  $C_1$ , policy  $S_1$  is proposed to reduce the answer time (the time period between the time when an order is requested and the time when the order is responded by the driver). This in turn meets more call orders requests. Both policy  $S_2$  in city  $C_2$  and policy  $S_3$  in city  $C_3$  are designed to guide drivers to regions with more orders in order to reduce drivers' idle time ratio. Policies  $S_2$  and  $S_3$  are designed to assign more drivers to areas with more orders. This in turn reduces drivers' downtime and increase their income. Policy  $S_4$  aims to balance drivers' downtime and their average pick-up distance.

We also apply our test to another four datasets collected from four A/A experiments which compare the standard policy against itself. These A/A experiments are conducted two weeks before the A/B experiments. Each lasts for 14 days and thirty-minutes is defined as one time unit. We remark that the A/A experiment is employed as a sanity check for the validity of the proposed test. We expect our test will not reject the null when applied to these datasets, since the sole standard policy is used.

We fit the proposed L-TVCDP models to each of the eight datasets. In Figures 7 and 8, we plot the predicted outcomes against the observed values and plot the corresponding residuals over time for policy  $S_1$ . Results for policies  $S_2$ – $S_4$  are represented in Figure S.5 in the supplementary article. It can be seen that the predicted outcomes are very close to the observed values, suggesting that the proposed model fits the data well. P-values of the proposed tests are reported in Tables 1 and 2. As expected, the proposed test does not reject the null hypothesis when applied to all datasets from A/A experiments. When applied to the data from A/B experiments, it can be seen that the new policy  $S_1$  directly improves the answer rate and the completion rate, while increasing drivers' total income in city  $C_1$ . It also significantly increases drivers' income in the long run. Policy  $S_2$  has significant direct and indirect effects on drivers' income as expected. Policy  $S_4$  significantly increases the immediate answer rate, while improving the overall passenger satisfaction. However, policy  $S_3$  is not significantly better than the standard policy.

We further apply the proposed test to two real datasets collected from an A/A and A/B experiment

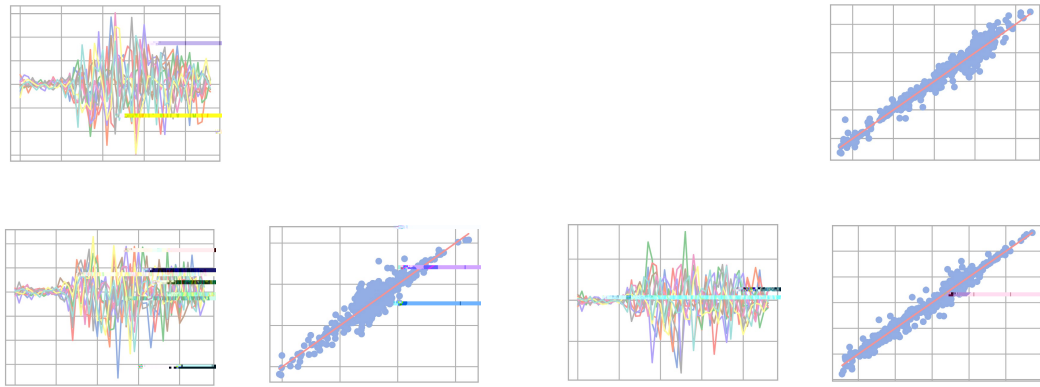
**Table 1.** One sided p-values of the proposed test for DE, when applied to eight datasets collected from the A/A or A/B experiment based on the temporal alternation design, with DTI, ART and CRT corresponding to drivers' total income, the answer rate and the completion rate, respectively.

	AA			AB	
	DTI(%)	ART(%)	CRT(%)	DTI(%)	ART(%)

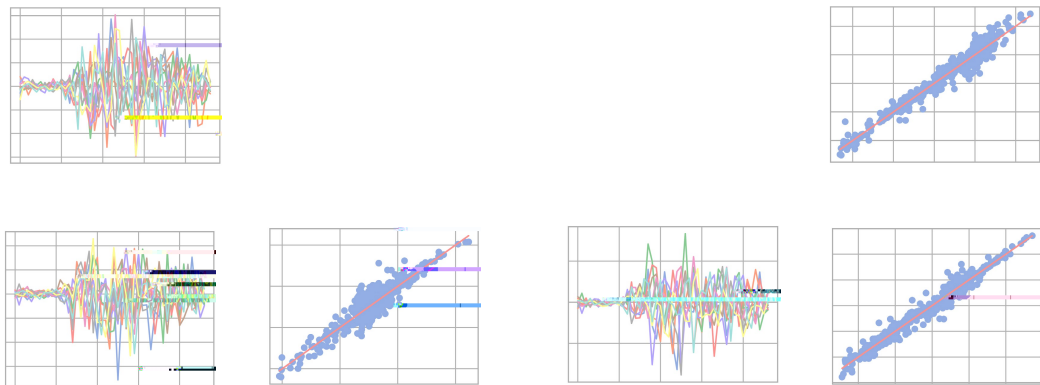
**Table 3.** One sided p-values of the proposed test, when applied to two datasets collected from the A/A or A/B experiment based on the spatio-temporal alternation design. Drivers’ total income is set to be the outcome of interest.

	DE		IE	
	AA	AB	AA	AB
p-value	0.176	0.001	0.334	0.000

under the spatio-temporal alternation design, conducted in city  $C_5$ . This city is partitioned into 17 regions. Within each region, more than 90% orders are answered by drivers in the same region. Similar to the temporal alternation design, both experiments last for 14 days and 30-minutes is set as one time unit. We take the number of requests as the state variables and drivers’ total income as the outcome, as in Section 5.2. In Figures 9 and 10, we plot the fitted drivers’ total income and the fitted number of requests against their observed values, and plot the corresponding residuals over time. We only present results associated with 2 regions in the city for space economy. The fitted values and residuals associated with other regions are similar and we do not present them to save space. It can be seen that the proposed models fit these datasets well. In addition, we report the p-values of the proposed test in Table 3. It can be seen that the new policy significantly increases drivers’ income. When applied to the dataset from the A/A experiment, it fails to reject either null hypothesis.



**Fig. 9.** Plots of the fitted drivers’ income against the observed values, as well as the corresponding residuals. Data are collected from an A/A or A/B experiment under the spatio-temporal alternation design.



**Fig. 10.** Plots of the fitted number of orders against the observed values, as well as the corresponding residuals. Data are collected from an A/A or A/B experiment under the spatio-temporal alternation design.

## 7. Discussion

In this study, driven by the need for policy evaluation in technological companies, we thoroughly examine AB testing for temporal and/or spatial dependent experiments, particularly in scenarios characterized by weak signals, (spatio)-temporal random effects, and intricate interference structures. Our approach offers two key benefits. Firstly, it accommodates the switchback design, which can significantly enhance testing power. As explained earlier, by applying diverse treatments to neighboring time points, we can potentially offset the impact of random effects at these times, resulting in more efficient estimations of treatment effects. Secondly, we break down the ATE into its DE and IE components. We then advocate for testing these effects separately. This separation aids decision-makers in gaining a clearer understanding of how different policies function and in devising more effective strategies and designs. Further details can be found in Section S.12.4 of the supplementary document.

There are several intriguing avenues for future research. Firstly, considering Assumptions 1 and 2, it's worth exploring scenarios where errors in the state regression model are not necessarily independent over time. This can be achieved by incorporating random effects into the state regression model, allowing for correlated errors over time. However, this introduces dependencies between these random effects, which in turn affects the conditional independence of past and future features. Consequently, the Markov assumption is violated, and applying existing OPE methods and our proposal from Section 2 directly would result in biased policy value estimations. In Section S.12.1 of the supplementary document, we present two approaches to mitigate this endogeneity bias. Secondly, we can delve into situations involving a large number of state variables. However, in ride-sharing platforms, it's reasonable to assume that the dimension of state variables is fixed. This typically involves a two-dimensional market feature, encompassing the number of call orders and the number of available drivers. We outline potential extensions to high-dimensional settings in Section S.12.2 of the supplementary document. Thirdly, while the interference structure examined in this work is general, it remains relatively simple. It would be intriguing to explore more complex structural interferences across both space and time. Lastly, addressing statistical inference for deep neural networks remains an open challenge. This could represent a significant step toward incorporating deep learning into causal inference, offering promising directions for future research.

## Acknowledgments

The first three authors, Drs. Luo, Yang, and Shi, contributed equally to this paper. This work was finished when Drs. Luo, Ye, and Zhu worked at Didi Chuxing. The authors would like to express their gratitude to the Editors, the Associate Editor and the two reviewers for their constructive and insightful comments.

## Supplementary material

Supplementary data is available online at Journal of the Royal Statistical Society online.

## Conflict of interests

None declared.

## Funding

Yang's research is partially supported by China Postdoctoral Science Foundation (No. 2022TQ0360, 2022M723334) and the National Natural Science Foundation of China (No. 71988101). Shi's research is partially supported by an EPSRC grant EP/W014971/1. The content is solely the responsibility of the authors and does not necessarily represent the official views of any funding agency.

## Data availability

The code can be found on our GitHub page at <https://github.com/BIG-S2/STVCM>. We have included detailed instructions in the "Readme" file on how to reproduce Figures 1-10 and Tables 1-3. The real data used in our study is proprietary and cannot be shared publicly. However, we have provided simulated datasets that can yield similar results for broader accessibility and understanding.

## References

- Alonso-Mora, J., Samaranayake, S., Wallar, A., Frazzoli, E. and Rus, D. (2017) On-demand high-capacity ride-sharing via dynamic trip-vehicle assignment. *Proceedings of the National Academy of Sciences*, **114**, 462–467.
- Arkhangelsky, D., Imbens, G. W., Lei, L. and Luo, X. (2021) Double-robust two-way-fixed-effects regression for panel data. *arXiv preprint arXiv:2107.13737*.
- Aronow, P. M. and Samii, C. (2017) Estimating average causal effects under general interference, with application to a social network experiment. *The Annals of Applied Statistics*, **11**, 1912–1947.
- Aronow, P. M., Samii, C. and Wang, Y. (2020) Design-based inference for spatial experiments with interference. *arXiv preprint arXiv:2010.13599*.
- Bakshy, E., Eckles, D. and Bernstein, M. S. (2014) Designing and deploying online field experiments. In *Proceedings of the 23rd International Conference on World Wide Web*, 283–292.
- Bimpikis, K., Candogan, O. and Saban, D. (2019) Spatial pricing in ride-sharing networks. *Operations Research*, **67**, 744–769.
- Bojinov, I. and Shephard, N. (2019) Time series experiments and causal estimands: exact randomization tests and trading. *Journal of the American Statistical Association*, **114**, 1665–1682.
- Boruvka, A., Almirall, D., Witkiewitz, K. and Murphy, S. A. (2018) Assessing time-varying causal effect moderation in mobile health. *Journal of the American Statistical Association*, **113**, 1112–1121.
- Castillo, J. C., Knoepfle, D. and Weyl, G. (2017) Surge pricing solves the wild goose chase. In *Proceedings of the 2017 ACM Conference on Economics and Computation*, 241–242.
- Cohen, M. C., Fiszler, M. D. and Kim, B. J. (2022) Frustration-based promotions: Field experiments in ride-sharing. *Management Science*, **68**, 2432–2464.
- De Chaisemartin, C. and d’Haultfoeuille, X. (2020) Two-way fixed effects estimators with heterogeneous treatment effects. *American Economic Review*, **110**, 2964–96.
- Garg, N. and Nazerzadeh, H. (2022) Driver surge pricing. *Management Science*, **68**, 3219–3235.
- Hagi, A. and Wright, J. (2019) The status of workers and platforms in the sharing economy. *Journal of Economics & Management Strategy*, **28**, 97–108.
- Halloran, M. E. and Hudgens, M. G. (2016) Dependent happenings: A recent methodological review. *Current Epidemiology Reports*, **3**, 297–305.
- Hu, Y., Li, S. and Wager, S. (2022) Average direct and indirect causal effects under interference. *Biometrika*, **109**, 1165–1172.
- Hu, Y. and Wager, S. (2021) Off-policy evaluation in partially observed markov decision processes under sequential ignorability. *arXiv preprint arXiv:2110.12343*.
- Hudgens, M. G. and Halloran, M. E. (2008) Toward causal inference with interference. *Journal of the American Statistical Association*, **103**, 832–842.
- Imai, K. and Kim, I. S. (2021) On the use of two-way fixed effects regression models for causal inference with panel data. *Political Analysis*, **29**, 405–415.
- Jiang, N. and Li, L. (2016) Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning*, 652–661. PMLR.
- Johari, R., Li, H., Liskovich, I. and Weintraub, G. Y. (2022) Experimental design in two-sided platforms: An analysis of bias. *Management Science*, **68**, 7065–7791.
- Kallus, N. and Uehara, M. (2020) Double reinforcement learning for efficient off-policy evaluation in markov decision processes. *Journal of Machine Learning Research*, **21**.
- (2022) Efficiently breaking the curse of horizon in off-policy evaluation with double reinforcement learning. *Operations Research*, **70**, 3282–3302.
- Lale, S., Azizzadenesheli, K., Hassibi, B. and Anandkumar, A. (2021) Adaptive control and regret minimization in linear quadratic gaussian (lqg) setting. In *2021 American Control Conference (ACC)*, 2517–2522. IEEE.
- Larsen, N., Stallrich, J., Sengupta, S., Deng, A., Kohavi, R. and Stevens, N. T. (2023) Statistical challenges in online controlled experiments: A review of a/b testing methodology. *The American Statistician*, 1–15.
- Lee, L. (2007) Identification and estimation of econometric models with group interactions, contextual factors and fixed effects. *Journal of Econometrics*, **140**, 333–374.
- Lewis, G. and Syrgkanis, V. (2020) Double/debiased machine learning for dynamic treatment effects via g-estimation. *arXiv preprint arXiv:2002.07285*.
- Liao, P., Klasnja, P. and Murphy, S. (2021) Off-policy estimation of long-term average outcomes with

- applications to mobile health. *Journal of the American Statistical Association*, **116**, 382–391.
- Liao, P., Qi, Z., Klasnja, P. and Murphy, S. (2020) Batch policy learning in average reward markov decision processes. *arXiv preprint arXiv:2007.11771*.
- Liu, L., Hudgens, M. G. and Becker-Dreps, S. (2016) On inverse probability-weighted estimators in the presence of interference. *Biometrika*, **103**, 829–842.
- Liu, Q., Li, L., Tang, Z. and Zhou, D. (2018) Breaking the curse of horizon: Infinite-horizon off-policy estimation. In *Advances in Neural Information Processing Systems*, vol. 31.
- Luckett, D. J., Laber, E. B., Kahkoska, A. R., Maahs, D. M., Mayer-Davis, E. and Kosorok, M. R. (2020) Estimating dynamic treatment regimes in mobile health using v-learning. *Journal of the American Statistical Association*, **115**, 692–706.
- Luedtke, A. R. and Van Der Laan, M. J. (2016) Statistical inference for the mean outcome under a possibly non-unique optimal treatment strategy. *The Annals of Statistics*, **44**, 713–742.
- Manski, C. F. (2013) Identification of treatment response with social interactions. *The Econometrics Journal*, **16**, S1–S23.
- Munro, E., Wager, S. and Xu, K. (2021) Treatment effects in market equilibrium. *arXiv preprint arXiv:2109.11647*.
- Perez-Heydrich, C., Hudgens, M. G., Halloran, M. E., Clemens, J. D., Ali, M. and Emch, M. E. (2014) Assessing effects of cholera vaccination in the presence of interference. *Biometrics*, **70**, 731–741.
- Pollmann, M. (2020) Causal inference for spatial treatments. *arXiv preprint arXiv:2011.00373*.
- Puelz, D., Basse, G., Feller, A. and Toulis, P. (2019) A graph-theoretic approach to randomization tests of causal effects under general interference. **arXiv**, 1910.10862v1.
- Puterman, M. L. (2014) *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.
- Qin, Z. T., Zhu, H. and Ye, J. (2022) Reinforcement learning for ridesharing: An extended survey. *Transportation Research Part C: Emerging Technologies*, **144**, 103852.
- Reich, B. J., Yang, S., Guan, Y., Giffin, A. B., Miller, M. J. and Rappold, A. (2020) A review of spatial causal inference methods for environmental and epidemiological applications. **arXiv**, 2007.02714v1.
- Rubin, D. (1980) Discussion of "randomization analysis of experimental data in the fisher randomization test" by d. basu. *Journal of the American Statistical Association*, **75**, 591–593.
- Sävje, F., Aronow, P. and Hudgens, M. (2021) Average treatment effects in the presence of unknown interference. *The Annals of Statistics*, **49**, 673–701.
- Schmidt-Hieber, J. (2020) Nonparametric regression using deep neural networks with relu activation function. *The Annals of Statistics*, **48**, 1875–1897.
- Shen, Z., Yang, H. and Zhang, S. (2019) Deep network approximation characterized by number of neurons. *arXiv preprint arXiv:1906.05497*.
- (2022) Optimal approximation rate of relu networks in terms of width and depth. *Journal de Mathématiques Pures et Appliquées*, **157**, 101–135.
- Shi, C., Wan, R., Chernozhukov, V. and Song, R. (2021) Deeply-debiased off-policy interval estimation. In *International conference on machine learning*, 9580–9591. PMLR.
- Shi, C., Wan, R., Song, G., Luo, S., Song, R. and Zhu, H. (2022a) A multi-agent reinforcement learning framework for off-policy evaluation in two-sided markets. *arXiv preprint arXiv:2202.10574*.
- Shi, C., Wang, X., Luo, S., Zhu, H., Ye, J. and Song, R. (2023) Dynamic causal effects evaluation in a/b testing with a reinforcement learning framework. *Journal of the American Statistical Association*, **118**, 2059–2071.
- Shi, C., Zhang, S., Lu, W. and Song, R. (2022b) Statistical inference of the value function for reinforcement learning in infinite-horizon settings. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, **84**, 765–793.
- Shumway, R. and Stoffer, D. (2010) *Time series analysis and its applications with R examples (3rd ed.)*. Springer.
- Sobel, M. E. (2006) What do randomized studies of housing mobility demonstrate?: Causal inference in the face of interference. *Journal of the American Statistical Association*, **101**, 1398–1407.
- 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**, 967–976.
- Tang, X., Qin, Z., Zhang, F., Wang, Z., Xu, Z., Ma, Y., Zhu, H. and Ye, J. (2019) A deep value-network based approach for multi-driver order dispatching. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 1780–1790.



- Tchetgen Tchetgen, E. J. and VanderWeele, T. J. (2012) On causal inference in the presence of interference. *Statistical Methods in Medical Research*, **21**, 55–75.
- Thomas, P. and Brunskill, E. (2016) Data-efficient off-policy policy evaluation for reinforcement learning. In *International Conference on Machine Learning*, 2139–2148. PMLR.
- Uehara, M., Shi, C. and Kallus, N. (2022) A review of off-policy evaluation in reinforcement learning. *arXiv preprint arXiv:2212.06355*.
- Van, D. and Wellner, J. A. (1996) *Weak convergence and empirical processes*. Springer,.
- Verbitsky-Savitz, N. and Raudenbush, S. W. (2012) Causal inference under interference in spatial settings: A case study evaluating community policing program in chicago. *Epidemiologic Methods*, **1**, 107–130.
- Wager, S. and Xu, K. (2021) Experimenting in equilibrium. *Management Science*, **67**, 6694–6715.
- Wooldridge, J. M. (2021) Two-way fixed effects, the two-way mundlak regression, and difference-in-differences estimators. *Available at SSRN 3906345*.
- Wu, C.-F. J. et al. (1986) Jackknife, bootstrap and other resampling methods in regression analysis. *The Annals of Statistics*, **14**, 1261–1295.
- Yan, S. and Yao, F. (2023) Nonparametric regression for repeated measurements with deep neural networks. *arXiv preprint arXiv:2302.13908*.
- Zhang, B., Tsiatis, A. A., Laber, E. B. and Davidian, M. (2013) Robust estimation of optimal dynamic treatment regimes for sequential treatment decisions. *Biometrika*, **100**, 681–694.
- Zhou, F., Luo, S., Qie, X., Ye, J. and Zhu, H. (2021) Graph-based equilibrium metrics for dynamic supply–demand systems with applications to ride-sourcing platforms. *Journal of the American Statistical Association*, **116**, 1688–1699.
- Zhu, H., Fan, J. and Kong, L. (2014) Spatially varying coefficient model for neuroimaging data with jump discontinuities. *Journal of the American Statistical Association*, **109**, 1084–1098.
- Zigler, C. M., Dominici, F. and Wang, Y. (2012) Estimating causal effects of air quality regulations using principal stratification for spatially correlated multivariate intermediate outcomes. *Biostatistics*, **13**, 289–302.