Optimal Experimental Design for Staggered Rollouts

Ruoxuan Xiong*

Department of Management Science & Engineering, Stanford University, rxiong@stanford.edu

Susan Athev

Graduate School of Business, Stanford University, athey@stanford.edu

Mohsen Bayati

Graduate School of Business, Stanford University, bayati@stanford.edu

Guido Imbens

Graduate School of Business, Stanford University, imbens@stanford.edu

Experimentation has become an increasingly prevalent tool for guiding decision-making and policy choices. A common hurdle in designing experiments is the lack of statistical power. In this paper, we study the optimal multi-period experimental design under the constraint that the treatment cannot be easily removed once implemented; for example, a government might implement a public health intervention in different geographies at different times, where the treatment cannot be easily removed due to practical constraints. The treatment design problem is to select which geographies (referred by units) to treat at which time, intending to test hypotheses about the effect of the treatment. When the potential outcome is a linear function of unit and time effects, and discrete observed/latent covariates, we provide an analytically feasible solution to the optimal treatment design problem where the variance of the treatment effect estimator is at most $1 + O\left(\frac{1}{N^2}\right)$ times the variance using the optimal treatment design, where N is the number of units. This solution assigns units in a staggered treatment adoption pattern – if the treatment only affects one period, the optimal fraction of treated units in each period increases linearly in time; if the treatment affects multiple periods, the optimal fraction increases non-linearly in time, smaller at the beginning and larger at the end. In the general setting where outcomes depend on latent covariates, we show that historical data can be utilized in designing experiments. We propose a data-driven local search algorithm to assign units to treatment times. We demonstrate that our approach improves upon benchmark experimental designs via synthetic interventions on the influenza occurrence rate and synthetic experiments on interventions for in-home medical services and grocery expenditure.

Key words: Experimental Design, Contamination and Spillover Effects, Staggered Rollouts, Treatment Effect Estimation, Stratification

1. Introduction

Experimentation is a cornerstone of decision-making. Large technology and retail companies run thousands or even tens of thousands of experiments per year to evaluate the impact of various decisions, design better products, and understand mechanisms (Muchnik et al. 2013, Aral and

^{*} Alphabetical order other than the first author.

Walker 2014, Panniello et al. 2014, Cheung et al. 2017, Brynjolfsson et al. 2018, Cohen et al. 2018, Brynjolfsson et al. 2019, Salganik 2019, Gordon et al. 2019, Joo et al. 2019, Zhang et al. 2019a,b, Ülkü et al. 2019, Lee and Hosanagar 2019, Cui et al. 2020, Sun et al. 2020a,b). In some cases, a variety of firm constraints lead firms (units) to introduce new products or features at different times. In some cases, a variety of practical constraints lead clinical researchers and policymakers to introduce the treatment or intervention at different times. After such a rollout has occurred, analysts treat this type of "staggered adoption" of the new feature as a "natural experiment" and attempt to estimate causal effects (Athey and Stern 2002, Athey and Imbens 2018, Athey et al. 2018). However, less research has been done to consider how to design such experiments in anticipation of precisely analyzing treatment effects.

In this paper, we study how to optimally design experiments in settings where the experimental designer can choose not just which units to treat but also when to treat them. Our goal is to minimize the variance of the estimate for the treatment effect by choosing the optimal multi-period treatment design. We maintain the assumption that during the course of the experiment, once treated, a unit continues being treated until the end of the experiment. This assumption is motivated by practical considerations such as difficulties communicating with experimental subjects or frictions in implementing changes; we also extend our results to the scenario where the treatment can be removed.

Our main motivation to study this question is the fundamental challenge in experimental design: statistical power (Pocock and Simon 1975, Cohen 1992, Kuhfeld 2005, Lawrie et al. 2015, Hemming et al. 2015, Li et al. 2018b, Bhat et al. 2019), which is the likelihood to detect the effect of a new decision and is also a concern in observational studies (Fogarty and Small 2016, Heng et al. 2019). Statistical power can be low due to various reasons, but the primary reason is a small sample size. For an essential class of questions where the treatment effect has spillover effects, the relevant unit of analysis is an aggregated set of samples, such as a metropolitan area for testing public health interventions to reduce the occurrence of an infectious disease such as flu. Another example is randomization at a clinic instead of a patient-level for a clinical trial. It compares different weightloss programs where patients may meet each other in the waiting room and communicate about their own strategies (Leahey et al. 2014). While randomizing at a higher level of aggregation can effectively avoid contamination and spillover effects, but it substantially reduces the sample size and statistical power.

Measuring outcomes for units over time allows the experimental designer to improve the statistical power, as repeated observations for the same units allow the analyst to learn key parameters of the data generating process, and thus reduce the error in estimating outcomes for units under counterfactual treatment status (Crowder and Hand 1990, Baltagi 2008, Hsiao 2014). However,

determining the optimal timing of treatment assignments is challenging for two reasons. First, it is difficult to find an explicit functional form for how the treatment design affects the variance of the treatment effect estimator (Hussey and Hughes 2007, Hemming et al. 2015, Li et al. 2018b). Second, assigning units to either the treated or control group at each point in time is an integer programming problem, which is, in general, computationally expensive to solve (Bertsimas et al. 2015, Kallus 2018, Bhat et al. 2019). For the first challenge, our analysis is based on the best linear unbiased estimator, which provides the explicit functional form and has the smallest variance among all linear unbiased estimators. The main contribution of this paper is to provide closed-form solutions for the second challenge. We provide theoretical guarantees for the optimality of our proposed multi-period experiment designs and support their practical relevance via extensive simulations on real data sets.

This paper's objective is to optimally allocate subjects from the control groups to the treated groups over multiple periods so that the variance of our treatment effect estimator is minimized. We restrict treated units to switch back to the control. We assume potential outcomes follow a two-way fixed-effect linear model, and possibly with additional observed or latent covariates. Furthermore, we assume the treatment effect is constant, and we want to estimate this object as precise as possible. This assumption is less restrictive than it looks like. When the treatment effect is heterogeneous and time-varying, an object of interest is the average treatment effect (Rosenbaum and Rubin 1983, Abadie et al. 2010). We can estimate the average effect through the constant treatment effect model by decomposing the treatment effect into two terms: 1. the average effect, which is constant; 2. the heterogeneous and time-varying part with zero mean, which can be incorporated into the error term. We analyze the optimal experimental design to minimize the variance of the best linear unbiased estimator (BLUE) for the treatment effect. We find a treatment design with a constant growth rate of the treated percentage is near-optimal. We propose an algorithm based on simulated annealing and the minimax decision criterion to actively improve this near-optimal solution using historical control data.

Specifically, our first main contribution is to provide simple sufficient conditions for the optimal treatment design when the potential outcome can have observed and latent covariates. In particular, we allow latent covariates to address the concern of model misspecification that we may not fully observe all the covariates that affect the potential outcome. Furthermore, if we have many observed covariates, and we do not know which matter, we could instead use a data-driven approach to estimate the latent covariates that can well summarize the information in the data. When covariates take finitely many values, in an optimal treatment design, the treated proportion (fraction of treated units in each period) should be equal for each cluster of units with equal covariate values (stratum). We can interpret this treatment design as stratified randomization in both

| Model Two-way fixed effects | | Two-way fixed effects with carryover treatment effects and discrete covariates |
|---|---|--|
| Optimal $\left\ \frac{2t-1}{2T} \right\ $ fraction treated for solution $\left\ \text{all } t \right\ $ in a T -period design | $\left \begin{array}{l} \frac{2t-1}{2T} \end{array} \right $ fraction treated for each stratum and all t | Nonlinear fraction treated in t with five stages for each stratum |

Table 1 This table summarizes the optimal treatment design in various potential outcome models. The base case has two-way (unit and time) fixed effects (Section 3.1). A more general case has two-way fixed effects with additional discrete observed and/or latent covariates (Section 3.2). The most general case allows the effect of a treatment carries over to multiple periods (Section 3.3).

unit and time dimensions, as a contrast to the conventional A/B testing that only randomizes in the unit dimension. If the treatment only affects *one* period, the optimal treated proportion increases linearly in time. However, if the effect of treatment carries over to future periods (Zhang et al. 2019a), the optimal treated proportion increases non-linearly in time, is convex at the beginning, and concave towards the end. When there is no feasible solution to allow the exact optimal treated proportion for each stratum, we provide a *rounding* scheme that yields a feasible solution whose estimation variance is within $1 + O\left(\frac{1}{N^2}\right)$ of the estimation variance of the optimal solution, where N is the number of units. Our main results are summarized in Table 1.

Our second main contribution is to propose a data-driven algorithm to efficiently find a better treatment design when covariates are latent or take infinitely many values. Our algorithm exploits covariate information in the historical control data and searches for a better treatment design based on simulated annealing and minimax decision rule, which is motivated by our theoretical analysis.

We illustrate our treatment design and algorithm's performance through synthetic experiments on three real data sets on studying interventions for reducing flu occurrence rate, interventions for reducing the utilization of home-visit medical services, and marketing strategies for promoting grocery store sales. The first example is motivated by several features of the critical problem of rigorously studying the effect of various interventions during a pandemic such as COVID-19. First, in such experiments, statistical power is low not only because we need to randomize at large geographical areas due to the contagious nature of the disease, but also because the treatment effect may be smaller than the (seasonal) fluctuations of the disease. Second, our proposed solutions in this paper allow experiments to start at the beginning of the peak season, continue during the peak season, and terminate at the end of the peak season. Third, the interventions usually have a carryover effect on the disease occurrence rate in multiple future periods. Overall, based on our theoretical results and empirical analysis of real data, we find that our proposed treatment designs substantially outperform benchmark designs.

The rest of the paper is organized as follows. Section 2 describes the model and our objective. The main results for the analysis of treatment designs are presented in Section 3, and Section 4 provides our algorithms for finding optimal treatment designs. Finally, in Section 5, we demonstrate the efficiency gain from our results and algorithms by analyzing synthetic experiments on interventions for reducing flu occurrence rate, using a large multi-year real data set. We defer synthetic experiments on a secondary data set for reducing utilization of home-visit services and a tertiary data set for promoting grocery store sales to Appendix F.

1.1. Related Literature

We mainly focus on the problem of multi-period experimental design when treated units cannot switch back to the control, which is closely related to the stepped wedge design in clinical trials (Brown and Lilford 2006). In the stepped wedge design, an intervention is sequentially rolled out to participants (more often to clusters than to individuals) over several periods. Hemming et al. (2015) point out that the stepped wedge clustered studies tend to have better statistical performance than parallel clustered studies when the intra-cluster correlations are larger. Optimal allocation of clusters in the stepped wedge design has been studied under the linear model with time fixed effects, random cluster effects (Hussey and Hughes 2007, Hemming et al. 2015), and additional random interaction effects (Li et al. 2018b). In comparison, our setting and results are more general in the following three facets: 1. We allow each unit to have its own (observed or latent) covariates; 2. We allow the treatment to have carryover effects; 3. We show how historical control data can be used to improve treatment design. The first generalization is universal in practice; indeed, the assumption in prior literature on the linearity of outcomes in time effects and cluster effects is restrictive. The second generalization is relevant in many applications, such as in clinical trials (Grizzle 1965, Wallenstein and Fisher 1977, Willan and Pater 1986). For the third generalization, historical control data is commonly available in clinical trials and public health, so it is natural to leverage that data to enhance experiment design.

Our problem also shares similarities with the literature on covariate balancing and matching in experimental design (Pocock and Simon 1975, Cook et al. 2002, Hu et al. 2012, Bertsimas et al. 2015, 2019, Kallus 2018, Bhat et al. 2019, Krieger et al. 2019) and in causal inference (Nikolaev et al. 2013, Imai and Ratkovic 2014, Sauppe et al. 2014, Sun and Nikolaev 2016, Fan et al. 2016, Sauppe and Jacobson 2017, Li et al. 2018a, Kallus et al. 2018, Zhao 2019, Kallus 2020). Our experimental design setup is different in that we allocate treatment for multiple units over multiple periods and the treatment cannot be removed once allocated. Compared with the literature on designing covariate balanced experiments (Bertsimas et al. 2015, 2019, Kallus 2018, Bhat et al. 2019), we allow covariates to be latent, and we propose to search for a better treatment design from historical control data that contains information about the latent covariates. Moreover, instead of directly solving the integer programming problem¹, we propose to start from the analytical solution

¹ This is, in general, computationally expensive. This can be infeasible when covariates are latent.

that is optimal in the model without covariates, and actively improve this solution via historical data in an efficient way.

Another relevant literature to our findings is prior work on stratification in randomized experiments, such as in clinical trials (Simon 1979, Polit and Beck 2008, Athey and Imbens 2017), and in stratified sampling (Mulvey 1983, Cheng and Davenport 1989, Fox 2000). They suggest the statistical power of the treatment effect estimator can be increased by partitioning the population into small strata, randomizing within each stratum, and adjusting the standard errors. Our results convey a similar insight; we prove that stratification increases power in multi-period experimental design with discrete covariates. Furthermore, for more general covariates, we propose a data-driven stratification heuristic based on historical data.

The rich literature in online learning and the multi-armed bandit (e.g., Bubeck and Cesa-Bianchi (2012), Lattimore and Szepesvári (2018)) can also be viewed as a form of multi-period experimental design. However, in that setting, the decisions (or treatment designs) are adaptively updated, given the experiment's partially observed results. In contrast, we consider settings where the experiment designer, should commit to a procedure for the future rounds, in advance.

2. Model and Problem Statement

In this section, we will introduce the potential outcome model under the control and treatment in Section 2.1. Next, in Section 2.2, we describe the main treatment allocation problem that needs to be solved to maximize the statistical power of the experiment.

2.1. Data Model

Before we massively scale up implementing a new policy, we want to evaluate its treatment effect. Assume we have N units to apply this new policy to, and we can observe their outcomes for T periods. We have the following assumption for the potential outcome.

Assumption 1 (Data Model Assumption). Assume the potential outcome for these N units over the T time periods can be modeled as

$$Y_{it}(z_{it}) = \alpha_i + \beta_t + L_{it} + \tau z_{it} + \varepsilon_{it} \tag{1}$$

and in matrix notation

$$Y(Z) = \alpha \vec{1}^{\top} + \vec{1}\beta^{\top} + L + \tau Z + \varepsilon, \qquad (2)$$

where τ is the treatment effect, $z_{it} \in \{-1,1\}$ indicates whether unit i at time t is treated $(z_{it} = 1)$ or not $(z_{it} = -1)$, α_i and β_t are unit and time fixed effects. Furthermore, we allow matrix $L = [L_{it}] \in \mathbb{R}^{N \times T}$ to be in one of the following two cases:

- 1. L = 0:
- 2. $L = X\theta^{\top} + UV^{\top}$, where $X \in \mathbb{R}^{N \times r}$ and $U = \begin{bmatrix} u_1^{\top} & \cdots & u_N^{\top} \end{bmatrix}^{\top} \in \mathbb{R}^{N \times k}$ are observed and latent covariates, $\theta \in \mathbb{R}^{T \times r}$ and $V = \begin{bmatrix} v_1^{\top} & \cdots & v_T^{\top} \end{bmatrix}^{\top} \in \mathbb{R}^{T \times k}$ are both unknown coefficients. We also assume that U is deterministic, and V is random and independent of the treatment with $\mathbb{E}[v_t | \alpha, \beta, X] = 0$ and $\text{Cov}[v_t | \alpha, \beta, X] = \Sigma_V$ for all t.

We assume L is low-rank, that is $r, k \ll \min(N, T)$. In the second case, r and k can be 0. Furthermore, we assume ε_{it} is the i.i.d. observed noise and independent of the treatment with $\mathbb{E}[\varepsilon_{it}|\alpha,\beta,X] = 0$ and $\text{Cov}[\varepsilon_{it}|\alpha,\beta,X] = \sigma^2$ for all i and t.

REMARK 1. The term UV^{\top} can be thought of as the structured component of the noise that models correlation in the noise for different units. Since in this structure the randomness is in V, the properties for U and V do not appear as symmetric variants of each other. Alternatively, we can assume V is deterministic, and U is random, which would model the case where noise has time-series correlations. We only focus on the former case since the latter would be similar. Another straightforward extension of our model that we do not cover in the paper is when there are two independent sources of structured noise (one with across unit correlations and one with across time correlations). Furthermore, the assumption that V has mean zero holds without loss of generality because otherwise, one can update the fixed effects α and β as well as U and V, so that this assumption holds. The assumption on the randomness of this model, V and ε , is standard and allows regression estimators to be unbiased that will be discussed in detail in Section 2.2.

In this paper, our primary focus is the second case $L = X\theta^{\top} + UV^{\top}$ where the potential outcome has covariates that can be either observed, latent, or both. However, we will discuss the simple case of L = 0 only to build intuition for the general case.

We assume the treatment effect τ is constant but Remark 2 below discusses how more general treatment effects can be also captured by this setup and Remark 3 below discusses the extension of model (1) to carryover treatment effects.

We explicitly specify fixed affects $\alpha \vec{1}^{\top}$ and $\vec{1}\beta^{\top}$ in model (2) because it is common that the potential outcomes across units and across time periods have different means and it can potentially reduce the estimation bias of parameters in model (2).

In this paper, we focus on the *irreversible* treatment adoption pattern, which is stated in the following assumption. However, for completeness, we also solve the case when this assumption fails (the treatment can be removed) in Appendix A.1.

Assumption 2 (Irreversible Treatment Adoption Pattern). We assume the treatment adoption pattern is such that once a unit adopts treatment, it stays treated afterwards, that is,

$$z_{it} \leq z_{i,t+1}$$
, for all i and $1 \leq t \leq T-1$

First, we allow units to progressively adopt treatment. Our setting is different from paired and stratified experiments, where the treatment status stays the same throughout the experiment. Our setting is similar to the crossover trial in a longitudinal study. However, we prohibit treated units from switching back to the control, which is closely related to the stepped wedge treatment design in randomized controlled trials (Brown and Lilford 2006, Hussey and Hughes 2007, Woertman et al. 2013, Hemming et al. 2015). This restriction is closely related to the staggered treatment adoption pattern in observational studies (Athey and Stern 1998, Athey et al. 2018, Athey and Imbens 2018).

We study the irreversible scenario mainly for two reasons: 1. There are some practical constraints restricting units from frequently switching between control and treatment. For example, policies at the cluster level, such as programs to improve case management skills of healthcare staff at the hospital level, can hardly be suspended or rolled back. 2. When the treated units switch back to control, they may not return to the original control status. For example, the improvement of case management skills of healthcare staff can be persistent even if the programs are suspended.

REMARK 2. The assumption that τ is constant can be relaxed as follows. If the treatment effect is heterogeneous and time-varying, denoted by τ_{it} , we can still estimate the average treatment effect τ (which is typically the main goal in observational studies) by writing $\tau_{it} = \tau + \delta_{it}$ where δ_{it} has mean 0, and incorporate δ_{it} into ε_{it} in model (2). This is because each entry of Z is ± 1 which means the new noise (instead of ε_{it}) would be $\varepsilon_{it} \pm \delta_{it}$. This would still be an iid noise model if δ_{it} has a symmetric distribution around the origin.

REMARK 3. In some cases, the treatment can affect multiple periods and we are interested in estimating carryover treatment effects. We can extend model (1) to $Y_{it}(z_{it}) = \alpha_i + \beta_t + L_{it} + \tau_1 z_{it} + \tau_2 z_{i,t-1} + \cdots + \tau_{\ell+1} z_{i,t-\ell} + \varepsilon_{it}$, where ℓ is the number of lagged periods to which the effect of the treatment can carry over. We also study the optimal experimental design for this carryover model in Section 3.3.

2.2. Objective

Our objective is to estimate the treatment effect τ as precisely as possible. We need to consider two aspects to achieve our objective

- 1. How should we estimate τ ?
- 2. How can we find the optimal design matrix $Z \in \mathbb{R}^{N \times T}$ that minimizes the variance of the treatment effect estimator $\hat{\tau}$?

In the second question, the variance comes from the random observational noise ε and the random V in model (2).

As a preparation to discuss the first question, we state the well-known Gauss-Markov Theorem.

LEMMA 1 (Gauss-Markov Theorem). Consider a linear model $\vec{y} = \mathbf{X}\beta + noise$ with $\mathbb{E}[noise|\mathbf{X}] = 0$ and $\operatorname{Cov}[noise|\mathbf{X}] = \Omega$. If $\Omega = \sigma^2 I$, then the ordinary least squares (OLS) estimator $\hat{\beta} = (\mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top}\vec{y}$ is the best linear unbiased estimator (BLUE). Otherwise, when Ω is not a multiple of the identity matrix, the generalized least squares (GLS) estimator $\hat{\beta} = (\mathbf{X}^{\top}\Omega^{-1}\mathbf{X})^{-1}\mathbf{X}^{\top}\Omega^{-1}\vec{y}$ is BLUE. Here, "best" means the estimator has the lowest variance among all unbiased linear estimators.

For the first question, in model (2), if either L=0 or $L=X\theta^{\top}$ (L only has observed covariates), then $\hat{\tau}$ estimated from OLS is BLUE. If $L=X\theta^{\top}+UV^{\top}$, we can use GLS to estimate τ from

$$Y(Z) = \alpha \vec{1}^{\top} + \vec{1}\beta^{\top} + X\theta^{\top} + \tau Z + e, \tag{3}$$

where we incorporate UV^{\top} into the "error" term e, i.e., $e = UV^{\top} + \varepsilon$. Lemma 1 together with Assumption 1 implies that $\hat{\tau}$ from GLS is BLUE. Hence, our theoretical analysis is based on OLS and GLS because they are BLUE.

In practice, \vec{e} 's covariance matrix Ω is usually unknown, where \vec{e} is the vectorization of matrix e in Eq. (3). We can use the feasible generalized least squares (feasible GLS), where we get an estimate of Ω using the residuals from OLS. When errors' covariance matrix can be consistently estimated, feasible GLS is asymptotically more efficient than OLS. However, the covariance matrix has the dimension $NT \times NT$ that can be very large, and feasible GLS can be less efficient in finite samples. To address the efficiency concern, we can impose some structure on Ω , such as the diagonal plus low-rank structure. Alternatively, we propose an estimator (motivated by low-rank matrix estimation literature in machine learning) that can be more efficient but may be biased. The details are presented in Sections 4 and 5.

For the second question, we can find the optimal treatment design by solving the following integer programming problem,

$$\min_{Z} \operatorname{Var}(\hat{\tau}) \tag{4}$$
s.t. $z_{it} \leq z_{i,t+1}$, for all i and $1 \leq t \leq T-1$

$$z_{it} \in \{-1,1\}, \text{ for all } i \text{ and } t.$$

Solving the integer program (4) is challenging for two reasons. First, we need to understand how $Var(\hat{\tau})$ changes with z_{it} . Second, the decision variable z_{it} is an integer and solving the integer program (4) can be computationally expensive. Specifically, even in the special case where T=1 and $\hat{\tau}$ is obtained via OLS, we will see later in the integer program (8) that integer program (4) reduces to the number partitioning problem (Hayes 2002, Mertens 2006) which is NP-hard (the numbers corresponding to Γ that can take infinitely many values).

We address the first challenge by analyzing the variance of $\hat{\tau}$ estimated from OLS or GLS. Specifically, this variance is a quadratic function of z_{it} . For the second challenge, we show in Section 3 that when $\hat{\tau}$ is estimated by OLS or GLS depending on the structure of L, we can provide a near-optimal analytical solution, which has a surprisingly simple form and significantly outperforms benchmark treatment designs in empirical applications as shown in Section 5.

Remark 4. In the remaining, for notation simplicity, " $z_{it} \leq z_{i,t+1}$ " stands for " $z_{it} \leq z_{i,t+1}$, for all i and $1 \leq t \leq T - 1$ ", and " $z_{it} \in \{-1,1\}$ " stands for " $z_{it} \in \{-1,1\}$, for all i and t".

3. Results

We provide sufficient conditions for the optimal solution of the integer program (4) under mild assumptions, when we use BLUE (either OLS or GLS) to estimate $\hat{\tau}$, in Section 3.1 and 3.2. The sufficient conditions require the treated proportion to grow *linearly* in time. Based on the sufficient conditions, we demonstrate how to find a feasible solution with objective value (variance) to be at most $1 + O\left(\frac{1}{N^2}\right)$ times the objective value of the global optimum.

Furthermore, we extend our analysis to the scenario where the treatment has carryover effects, and our objective is to minimize the variance of estimated direct and lagged treatment effects. We find that the optimal proportion is *non-linear* in time. The details are presented in Section 3.3.

3.1. Optimal Solutions for Models without Covariates

In this section, we study the first case, i.e., L=0 (no observed or latent covariates). Model (1) can be rewritten as

$$Y_{it}(z_{it}) = \alpha_i + \beta_t + \tau z_{it} + \varepsilon_{it},$$

$$= \underbrace{\left[z_{it} \ \Gamma_{it}^{\top}\right]}_{1 \times (N+T)} \underbrace{\left[\begin{matrix} \tau \\ \eta \end{matrix}\right]}_{(N+T) \times 1} + \varepsilon_{it}$$
(5)

In matrix notation, we have

$$\vec{y} = \begin{bmatrix} \vec{z} \ \Gamma \end{bmatrix} \begin{bmatrix} \tau \\ \eta \end{bmatrix} + \vec{\varepsilon}, \tag{6}$$

where $\vec{y} = \begin{bmatrix} \vec{y}_1^\top \ \vec{y}_2^\top \cdots \vec{y}_T^\top \end{bmatrix}^\top \in \mathbb{R}^{NT \times 1}, \ \vec{y}_t = \begin{bmatrix} Y_{1t} \cdots Y_{Nt} \end{bmatrix}^\top \in \mathbb{R}^N, \ \tilde{I} = \begin{bmatrix} I_{N-1} \ \vec{0} \end{bmatrix}^\top \in \mathbb{R}^{N \times (N-1)}, \ \vec{0} = \{0\}^{N-1}, \ \vec{1} = \{1\}^N,$

$$\eta = \begin{bmatrix} \alpha \\ \beta \end{bmatrix}, \quad \alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_{N-1} \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_T \end{bmatrix}, \quad \text{and } \Gamma = \begin{bmatrix} \Gamma_{:1} \\ \Gamma_{:2} \\ \vdots \\ \Gamma_{:T} \end{bmatrix} = \begin{bmatrix} \tilde{I} & \tilde{I} \\ \tilde{I} & \tilde{I} \\ \vdots & \ddots \\ \tilde{I} & & \tilde{I} \end{bmatrix}.$$

Note that we restrict $\alpha_N = 0$ such that all other α_i and β_t can be uniquely identified. We use OLS to estimate τ and η . With some algebra, we have the following expression for the variance of $\hat{\tau}$,

$$\operatorname{Var}(\hat{\tau}) = \frac{\sigma^2}{\vec{z}^{\top} (I_{NT} - \Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top})\vec{z}}.$$
 (7)

Details to show Eq. (7) are presented in the Appendix C. With Eq. (7), the integer program (4) is simplified to

$$\min_{\vec{z}=[z_{it}]} \vec{z}^{\top} \Gamma(\Gamma^{\top} \Gamma)^{-1} \Gamma^{\top} \vec{z}
\text{s.t.} \quad z_{it} \leq z_{i,t+1}
z_{it} \in \{-1,1\}$$
(8)

Now we formally define our convex relaxation of the integer program (8). Let $\omega_t = \frac{1}{N} \sum_{i=1}^{N} z_{it}$ for all t. Therefore, it is easy to see that ω_t takes value from $\{-1, -1 + \frac{2}{N}, \dots, 1 - \frac{2}{N}, 1\}$, given the constraints of the integer program (8). We then relax the range of ω_t to $-1 \le \omega_t \le 1$ and obtain the following quadratic programming problem

$$\min_{\vec{\omega} \in \mathbb{R}^T} \vec{\omega}^{\mathsf{T}} P_{\vec{1}} \vec{\omega} + 2 \vec{d}^{\mathsf{T}} \vec{\omega}
\text{s.t.} \quad -1 \le \omega_t \le 1
\omega_t < \omega_{t+1},$$
(9)

with $P_{\vec{1}} = I - \frac{1}{T} \vec{1} \vec{1}^{\top} \in \mathbb{R}^{T \times T}$, $\vec{d} = [d_t] \in \mathbb{R}^T$, and $d_t = \frac{T+1-2t}{T}$ for all t. The following theorem provides the optimal treatment design for the quadratic program (9).

THEOREM 1 (Fixed Effects Only). Under Assumptions 1-2, and when L = 0, then an optimal solution for the quadratic program (9) is $\vec{\omega}^* \in \mathbb{R}^T$ that is defined by

$$\omega_t^* = \frac{2t - 1 - T}{T}, \quad \text{for all } t.$$
 (10)

Theorem 1 states that the analytical optimal solution ω_t^* of the quadratic program (9) is *linear* in t. Given ω_t^* , the optimal treated proportion is $\frac{1+\omega_t^*}{2}$. Figure 1 visualizes the optimal treated proportion at every time period for a T-period treatment design problem, where T ranges from 2 to 20. The treated proportion increases at the *constant rate* $\frac{1}{T}$.

Intuition. In order to build intuition for Theorem 1, the following two lemmas provide the optimal solution for two special cases of the problem.

LEMMA 2 (Time Fixed Effects Only). Consider the same assumptions as in Theorem 1, and in addition, assume that $Y_{it}(z_{it}) = \beta_t + \tau z_{it} + \varepsilon_{it}$. Then, any treatment design is optimal if it satisfies $\omega_t = 0$ for all t and $z_{it} \leq z_{i,t+1}$ for all i and t.

Lemma 2 states that if there are only time fixed effects in the potential outcome, the optimal treatment design is the parallel treatment design that half of the units are assigned treatment while the others are functioning as the control for all time periods. The intuition for Lemma 2 is that given the time effect β_t , observations within in a period serve as their own control to estimate τ . Formally, the variance only depends on the cross-sectional average of z_{it} . The variance is minimized when there are 50% treated and 50% control at every time period. Rows of Z are exchangeable, so only the treated proportion rather than whom to treat matters. Table 2a provides an example of the optimal treatment design. Lemma 2 implies that it is optimal to randomly and evenly split units into the treated and control groups from the beginning as in the conventional A/B testing and placebo tests.

LEMMA 3 (Unit Fixed Effects Only). Consider the same assumptions as in Theorem 1, and in addition, assume that $Y_{it}(z_{it}) = \alpha_i + \tau z_{it} + \varepsilon_{it}$. Then, any treatment design is optimal if it satisfies $\omega_t = -1$ for all $t < \frac{T+1}{2}$ and $\omega_t = 1$ for all $t > \frac{T+1}{2}$.

Lemma 3 states that if there are only unit fixed effects, the optimal treatment design is to allocate treatment to all units at halftime. The intuition for Lemma 3 is symmetric to that for Lemma 2: given the unit effect α_i , every unit's observations serve as their own control to estimate τ . Formally, the variance only depends on the *time-series average of* z_{it} . The variance is minimized when there are 50% treated and 50% control for each unit. Table 2b provides an example of the optimal treatment design.

Intuitively, when we have both unit and time fixed effects, the optimal treatment design is in between the optimal treatment designs in Lemmas 2 and 3. Specifically, the variance should depend on the average of z_{it} over time and the average of z_{it} over units. Under the irreversible treatment adoption pattern, the average of z_{it} over time can be identified, given the average of z_{it} over units. Therefore we can cancel out the average of z_{it} over time. The Hessian of the objective function in the quadratic program (9) is the positive semidefinite matrix $P_{\vec{1}}$, so the objective function is convex.

The analytical optimal solution ω_t^* is linear in t, and balances the learning of both unit and time fixed effects. If we aggregate observations for all time periods, there are 50% control and 50% treated. Table 2c provides an example of the optimal treatment design.

The optimal solution in Theorem 1, $\omega_t^* = \frac{2t-1-T}{T}$, is equivalent to $\frac{N(2t-1)}{2T}$ treated units at time period t. However, $\frac{N(2t-1)}{2T}$ may not be an integer. In order to get a feasible solution of the integer program (4), we use the following nearest integer rounding rule.

² If T is odd, units can be either treated or untreated at time period $t = \frac{T+1}{2}$.

| -1 -1 | -1 1 | -1 -1 |
|------------------------|------------------------|---------------------------|
| -1 -1 | -1 1 | -1 -1 |
| -1 -1 | -1 1 | -1 1 |
| -1 -1 | -1 1 | -1 1 |
| 1 1 | -1 1 | -1 1 |
| 1 1 | -1 1 | -1 1 |
| 1 1 | -1 1 | 1 1 |
| 1 1 | -1 1 | 1 1 |
| (a) Time fixed effects | (b) Unit fixed effects | (c) Two-way fixed effects |

Table 2 An optimal treatment design with N = 8 and T = 2 in the models with time fixed effects only, unit fixed effects only, and two-way fixed effects. Each row denotes a unit and each column denotes a time period.

Nearest integer rounding rule: We obtain a feasible solution Z^{rnd} by rounding $\frac{N(2t-1)}{2T}$ to the nearest integer. Specifically, $Z^{\text{rnd}} = [z_{it}^{\text{rnd}}]$ is a treatment design satisfying the constraints in the integer program (4) and

$$\frac{1}{N} \sum_{i=1}^{N} z_{it}^{\text{rnd}} = \begin{cases} \lfloor \frac{N(2t-1)}{2T} \rfloor & \text{if } \operatorname{frac}\left(\frac{N(2t-1)}{2T}\right) < \frac{1}{2}, \text{ or } \operatorname{frac}\left(\frac{N(2t-1)}{2T}\right) = \frac{1}{2} \text{ with } \frac{2t-1}{2T} < \frac{1}{2} \\ \lceil \frac{N(2t-1)}{2T} \rceil & \text{if } \operatorname{frac}\left(\frac{N(2t-1)}{2T}\right) > \frac{1}{2}, \text{ or } \operatorname{frac}\left(\frac{N(2t-1)}{2T}\right) = \frac{1}{2} \text{ with } \frac{2t-1}{2T} \ge \frac{1}{2}, \end{cases}$$

where $\operatorname{frac}\left(\frac{N(2t-1)}{2T}\right) = \frac{N(2t-1)}{2T} - \lfloor \frac{N(2t-1)}{2T} \rfloor$. The following theorem provides an upper bound for $\operatorname{Var}(\hat{\tau})$, when treatment design Z^{rnd} is used.

THEOREM 2. Under Assumption 1, let Z^* be the optimal solution of the integer program (4) and Z^{rnd} be the feasible solution of the integer program (4) using the nearest integer rounding rule. Denote the variance of $\hat{\tau}$ using the treatment design Z as $\text{Var}_Z(\hat{\tau})$. We have

$$\operatorname{Var}_{Z^{\operatorname{rnd}}}(\hat{\tau}) \leq \frac{1}{1 - 1/N^2} \operatorname{Var}_{Z^*}(\hat{\tau}).$$

Theorem 2 states that the nearest integer rounding rule outputs a solution that approximates the global optimum within a factor of $1 + O\left(\frac{1}{N^2}\right)$.

3.2. Optimal Solutions for Models with Covariates

In this section, we study the second case, i.e., $L = X\theta^{\top} + UV^{\top}$. Model (1) can be rewritten as

$$Y_{it}(z_{it}) = \alpha_i + \beta_t + X_i \theta_t + \tau z_{it} + \underbrace{u_i^\top v_t + \varepsilon_{it}}_{e_{it}}$$
(11)

Since e_{it} has a factor structure, e_{it} are correlated in the unit dimension. We use GLS to estimate τ , because GLS is BLUE when errors are correlated. Then we have the following theorem that provides the sufficient conditions and optimal treatment design for model (11).

THEOREM 3 (Observed and Latent Covariates). Suppose Assumptions 1-2 hold, rows in X have $\sum_{i=1}^{N} X_i = \vec{0}$, rows in U have $\sum_{i=1}^{N} u_i = \vec{0}$ and $\sum_{i=1}^{N} X_i u_i^{\top} = \mathbf{0}_{r,k}$, v_t 's covariance $\Sigma_V = I_k$,

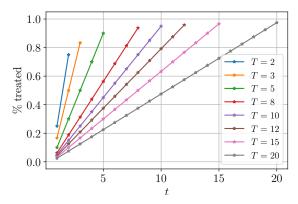


Figure 1 Optimal treated proportion at every time period for a T-period treatment design problem. The legends with different colors denote the number of periods T. For example, the leftmost blue line represents the optimal treated proportion at time periods 1 and 2 for a two-period treatment design problem, while the rightmost grey line represents the optimal treated proportion at every time period for a 20-period treatment design problem.

 $\operatorname{Cov}[v_t, v_{t-q} | \alpha, \beta, X] = 0$ for any q, $\operatorname{Cov}[v_t, \varepsilon_{is} | \alpha, \beta, X] = 0$ for any s, t and i, and ε_{it} 's variance $\sigma^2 = 1$. Then, Z is an optimal treatment design if it satisfies

$$\frac{1}{N} \sum_{i=1}^{N} z_{it} = \frac{2t - 1 - T}{T}, \qquad \frac{1}{N} \sum_{i=1}^{N} X_i z_{it} = \mu_X, \quad \frac{1}{N} \sum_{i=1}^{N} u_i z_{it} = \mu_U, \quad \text{for all } t$$
 (12)

for some $\mu_X \in \mathbb{R}^r$ and $\mu_U \in \mathbb{R}^k$, and if it satisfies the constraints in the integer program (4). Moreover, if (X_i, u_i) can only take G values, denoted as $(x_1, u_{0,1}), \dots, (x_G, u_{0,G})$, each with positive probability, we define $\omega_{g,t} = \frac{1}{|\mathcal{O}_g|} \sum_{i \in \mathcal{O}_g} z_{it}$, where $\mathcal{O}_g = \{i : (X_i, u_i) = (x_g, u_{0,g})\}$. Then any treatment design Z is optimal if it satisfies the constraints in the integer program (4) and

$$\omega_{g,t} = \frac{2t - 1 - T}{T}, \text{ for all } t \text{ and } g.$$
(13)

REMARK 5. The dimensions of X_i and u_i (i.e., r and k) can be zero. They correspond to the setting that the model only has latent covariates (r=0) and only has observed covarates (k=0). $\mathbf{0}_{r,k}$ is a zero matrix with dimension $r \times k$.

REMARK 6. When there are latent covariates u_i in the model, we need to know their value to stratify. If we have no information about u_i , we could only stratify over X_i and each stratum has treated proportion $\frac{2t-1}{2T}$ (in this context, units in a stratum have the same X_i , but may have different u_i). In many applications, we have historical control data, which contains information about u_i . We propose an algorithm in Section 4 to mimic stratification using historical data.

The assumptions in Theorem 3 can hold without loss of generality. First, for the assumption $\sum_{i=1}^{N} X_i = \vec{0}$ (X is orthogonal to $\vec{1}$) to hold, we can always use the Gram-Schmidt procedure³ to

³ Equivalent to QR decomposition

decompose $[\vec{1} \ X] = QR \in \mathbb{R}^{N \times (r+1)}$. Then the first column in Q is $\vec{1}$ and the second to last columns in Q are orthogonal to $\vec{1}$, which can be used as the observed covariates in model (11). We can use a similar procedure to project out the cross-section mean of u_i and $X_i u_i^{\top}$ for the assumptions $\sum_{i=1}^N u_i = \vec{0}$ and $\sum_{i=1}^N X_i u_i^{\top} = \mathbf{0}_{r,k}$ to hold. Second, if v_t has mean 0 and variance Σ_V for any positive definite matrix Σ_V , we can right multiply u_i^{\top} by $\Sigma_V^{-1/2}$ and left multiply v_t by $\Sigma_V^{-1/2}$ so that v_t has variance I_k (conditions $\sum_{i=1}^N u_i = \vec{0}$ and $\sum_{i=1}^N X_i u_i^{\top} = \mathbf{0}_{r,k}$ stay valid after this manipulation). Note that the assumptions on U and V are asymmetric and we provide a discussion in Remark 1. Moreover, we can multiply both the LHS and RHS of model (11) by $1/\sigma$ so that the variance of ε_{it} is 1 and the optimal treatment design does not change by this scaling. Moreover, we assume u_i is deterministic and v_t is random, following the literature on large dimensional factor modeling (Bai and Ng 2002, Bai 2003). We provide additional discussion of the model and assumptions in Theorem 3 in Appendix A.2.

Theorem 3 states that when covariates X_i and u_i only take finitely many values, we provide an optimal treatment design. This treatment design has linear staggered rollouts with stratification. In other words, for each stratum, the group of units with the same covariate value, we randomize the treatment allocation over time such that the treated proportion satisfies $\frac{2t-1}{2T}$ at time t. If we aggregate all strata, the overall treated proportion satisfies $\frac{2t-1}{2T}$ at time t, which is the same as Theorem 1. Table 3 provides an example of the optimal treatment design.

3.2.1. Overview of the Proof of Theorem 3

The complete proof is deferred to Appendix D and here we provide a high-level overview of the main steps. First, via algebraic manipulations, we show that $\operatorname{Var}(\hat{\tau}) = \vec{z}^{\top} (\Sigma_e^{-1} - \Sigma_e^{-1} \Gamma(\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1}) \vec{z}$ where Γ is the more general (with covariates) version of Γ from Section 3.1 and Σ_e is the covariance matrix of the vectorization of $e.^4$ Second, we show $\operatorname{min} \operatorname{Var}(\hat{\tau})$ can be simplified to minimize

$$\underbrace{(\vec{\omega}^{(1)})^{\top} P_{\vec{1}} \vec{\omega}^{(1)} + 2 \vec{d}^{\top} \vec{\omega}^{(1)}}_{(a)} + \underbrace{\sum_{j=2}^{r+1} (\vec{\omega}^{(j)})^{\top} P_{\vec{1}} \vec{\omega}^{(j)}}_{(b)} + \underbrace{\vec{z}^{\top} M_{U} \vec{z}}_{(c)}$$

where $\vec{\omega}_t^{(1)} = \frac{1}{N} \sum_{i=1}^N z_{it}$, $\vec{\omega}_t^{(j+1)} = \frac{1}{N} \sum_{i=1}^N X_{ij} z_{it}$, $P_{\vec{1}} = I_T - \frac{1}{T} \vec{1} \vec{1}^{\top}$, $d_t = \frac{T+1-2t}{T}$ and $M_U = P_{\vec{1}} \otimes U(I_k + U^{\top}U)^{-1}U^{\top}$. If we can find a solution that separately minimizes (a), (b) and (c), then this solution minimizes their sum. Third, (a) is minimized when Z satisfies $\frac{1}{N} \sum_{i=1}^N z_{it} = \frac{2t-1-T}{T}$ following Theorem 1. (b) depends on the cross-sectional average of z_{it} weighted by X_i and is minimized when Z satisfies $\frac{1}{N} \sum_{i=1}^N X_i z_{it} = \mu_X$, equivalently $\vec{\omega}^{(j+1)} = \mu_{X,j} \cdot \vec{1}$, where $\mu_{X,j}$ is the j-th entry in μ_X . (c) is

⁴ In practice, if we do not know Σ_e , we can first run OLS and then use the residual to estimate Σ_e , which is a consistent estimator for Σ_e .

(a) Two-way fixed effects with observed covariates X_i (b) Two-way fixed effects with latent covariates u_i

Table 3 An optimal treatment design with N=8 and T=2 under the model with either observed or latent covariates. In this toy example, there are two stratum, X_i (or u_i) = 1 and X_i (or u_i) = -1.

minimized when Z satisfies $\frac{1}{N} \sum_{i=1}^{N} u_i z_{it} = \mu_U$. Fourth, if (X_i, u_i) can only take G different values and there is a Z that satisfies $\omega_{g,t} = \frac{2t-1-T}{T}$ for all g and t, we plug this Z into the LHS of the sufficient conditions (12) and verify these conditions are satisfied.

3.2.2. Integrality Gap

In practice, we may not be able to find a feasible Z that satisfies Eq. (13) because $\frac{|\mathcal{O}_g|(2t-1)}{2T}$ may not be an integer for each strata, and therefore there does not exist a Z that satisfies Eq. (13). In the case where we only have X_i (no u_i), we use the same nearest integer rounding rule as the one in Section 3.1 to get a feasible treatment design Z^{rnd} of the integer program (4). Then we have the following theorem to provide a guarantee for Z^{rnd} .

THEOREM 4. Under the assumptions in Theorem 3 and $L = X\theta^{\top}$, let Z^* be the optimal solution of the integer program (4) and $Z^{\rm rnd}$ be a feasible solution of the integer program (4) using the nearest integer rounding rule. Denote the variance of $\hat{\tau}$ using the treatment design Z as ${\rm Var}_Z(\hat{\tau})$, $x_{j,\max} = \max_g |x_{gj}|$ and $N_{\min} = \min_g |\mathcal{O}_g|$. Assume $(1 + \sum_{j=1}^r x_{j,\max}^2)/N_{\min}^2 < 1$, we have

$$\operatorname{Var}_{Z^{\operatorname{rnd}}}(\hat{\tau}) \leq \frac{1}{1 - (1 + \sum_{j=1}^{r} x_{j,\max}^2) / N_{\min}^2} \operatorname{Var}_{Z^*}(\hat{\tau}).$$

Theorem 4 states the nearest integer rounding rule outputs a feasible solution whose objective value is at most $1+O\left(\frac{1}{N_{\min}^2}\right)$ times the global optimum when r and $x_{j,\max}$ are bounded. Since each realization of X_i takes positive probability bounded away from 0, we have $1+O\left(\frac{1}{N_{\min}^2}\right)=1+O\left(\frac{1}{N^2}\right)$ when $N\to\infty$. However, in finite samples, if G is large compared with N, the feasible solution may not be close to the global optimum. Then we could use partition methods, such as K-means, to partition units into G' groups based on their covariates' values. Within each group, we can randomize the treatment allocation such that Eq. (13) is satisfied. On the other hand, if X_i takes infinitely many values, we can also use this method.

3.3. Extension to Carryover Treatment Effects

In some cases, the effect of a treatment carries over to multiple periods, for example, in medical and clinical research. Researchers often need to consider the existence of carryover effects when designing cross-over trials that are widely used in medical and clinical research (Grizzle 1965, Wallenstein and Fisher 1977, Hills and Armitage 1979, Willan and Pater 1986, Freeman 1989, Senn 2002, Jones and Kenward 2014). The advertising effect on sales usually has a duration (Assmus et al. 1984), typically between six and nine months as shown by Leone (1995). In this section, we generalize model (1) and allow the treatment to have carryover effects,

$$Y_{it}(z_{it}) = \alpha_i + \beta_t + X_i \theta_t + \tau_1 z_{it} + \tau_2 z_{i,t-1} + \dots + \tau_{\ell+1} z_{i,t-\ell} + u_i^{\top} v_t + \varepsilon_{it}, \tag{14}$$

where $\varepsilon_{it} \stackrel{i.i.d.}{\sim} (0, \sigma^2)$. In this model, we use GLS to estimate $\tau_1, \dots, \tau_{\ell+1}$, and our objective is to minimize the estimation variance of $\tau_1, \dots, \tau_{\ell+1}$ by finding the optimal treatment design Z. We stack the observations Y_{it} from time $\ell+1$ to T together and in matrix notation

$$\underbrace{\begin{bmatrix} \vec{y}_{\ell+1} \\ \vec{y}_{\ell+2} \\ \vdots \\ \vec{y}_T \end{bmatrix}}_{\vec{\eta}} = \underbrace{\begin{bmatrix} \vec{z}_1 & \vec{z}_2 & \cdots & \vec{z}_{\ell+1} \\ \vec{z}_2 & \vec{z}_3 & \cdots & \vec{z}_{\ell+2} \\ \vdots & \vdots & \ddots & \vdots \\ \vec{z}_{T-\ell} & \vec{z}_{T-\ell+1} & \cdots & \vec{z}_T \end{bmatrix}}_{\vec{z}} \underbrace{\begin{bmatrix} \tau_{\ell+1} \\ \vdots \\ \tau_1 \end{bmatrix}}_{\vec{\tau}} + \Gamma \begin{bmatrix} \alpha \\ \tilde{\beta} \end{bmatrix} + \underbrace{\begin{bmatrix} \vec{e}_{\ell+1} \\ \vec{e}_{\ell+2} \\ \vdots \\ \vec{e}_T \end{bmatrix}}_{\vec{r}}, \text{ where } \vec{z} = \begin{bmatrix} z_{1t} \\ z_{2t} \\ \vdots \\ z_{Nt} \end{bmatrix}, \ \tilde{\beta} = \begin{bmatrix} \tilde{\beta}_{\ell+1} \\ \vdots \\ \tilde{\beta}_T \end{bmatrix}, \ \vec{e}_t = \begin{bmatrix} e_{1t} \\ e_{2t} \\ \vdots \\ e_{Nt} \end{bmatrix},$$

 $\tilde{\beta}_t = \begin{bmatrix} \beta_t \ \theta_t^{\top} \end{bmatrix}^{\top}, \vec{y}_t$ is the same as the \vec{y}_t in Eq. (6), and $\Gamma \in \mathbb{R}^{(N(T-\ell)) \times (N+(T-\ell-1)p)}$ is the more general (with covariates) version of the Γ in Eq. (6).

Our objective is to minimize $Var(\hat{\vec{\tau}})$, where

$$\operatorname{Var}(\hat{\vec{\tau}}) = (\mathcal{Z}^{\top} \left(\Sigma_e^{-1} - \Sigma_e^{-1} \Gamma (\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1} \right) \mathcal{Z})^{-1},$$

 Σ_e is the covariance matrix of the vectorization of $e = [e_{it}]$ and $e_{it} = u_i^{\top} v_t + \varepsilon_{it}$. Note that $\operatorname{Var}(\hat{\vec{\tau}})$ is a matrix and minimizing a matrix is not a well-defined problem. Instead, we consider minimizing $\operatorname{tr}(\operatorname{Var}(\hat{\vec{\tau}}))$ because the objects of interest are $\operatorname{Var}(\hat{\tau}_l)$ (the diagonal entries in $\operatorname{Var}(\hat{\vec{\tau}})$) and minimizing $\operatorname{tr}(\operatorname{Var}(\hat{\vec{\tau}}))$ is a well-defined problem. A surrogate for minimizing $\operatorname{tr}(\operatorname{Var}(\hat{\vec{\tau}}))$ is maximizing $\operatorname{tr}(\operatorname{Var}(\hat{\vec{\tau}}))$, that is,

$$\max_{Z} \operatorname{tr}(\mathcal{Z}^{\top} \left(\Sigma_{e}^{-1} - \Sigma_{e}^{-1} \Gamma(\Gamma^{\top} \Sigma_{e}^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_{e}^{-1} \right) \mathcal{Z})$$
s.t. $z_{it} \leq z_{i,t+1}$

$$z_{it} \in \{-1, 1\}$$

The optimal solution Z of the integer program (24) is called the $\mathbf{T}(\text{trace})$ -optimal treatment design. We also consider another objective, minimizing $\det(\text{Var}(\hat{\vec{\tau}}))$, in Appendix A.3. The following theorem provides sufficient conditions for the optimal solution of the integer program (24) as well as the analytical optimal solution when covariates can only take finitely many values.

THEOREM 5 (Carryover Effects with Observed and Latent Covariates). Suppose the potential outcome follows model (14), assumptions in Theorems 3 hold and $T > \frac{\ell^3 + 13\ell^2 + 7\ell + 3}{8\ell}$. Z is an optimal treatment design of the integer program (15) if it satisfies

$$\frac{1}{N} \sum_{i=1}^{N} z_{it} = \omega_t^*, \qquad \frac{1}{N} \sum_{i=1}^{N} X_i z_{it} = \mu_X, \quad \frac{1}{N} \sum_{i=1}^{N} u_i z_{it} = \mu_U, \quad \text{for all } t$$
 (16)

for some $\mu_X \in \mathbb{R}^r$ and $\mu_U \in \mathbb{R}^k$ and the constraints in the integer program (4), where $\vec{\omega}^* \in \mathbb{R}^T$ is defined by

$$\omega_{t}^{*} = \begin{cases} -1 & t \leq \lfloor \ell/2 \rfloor \\ ((A^{(\ell)})^{-1}b^{(\ell)})_{t-\lfloor \ell/2 \rfloor} & \lfloor \ell/2 \rfloor < t \leq \ell \\ -1 + \frac{2t - (\ell+1)}{T - \ell} & \ell < t \leq T - \ell \\ -((A^{(\ell)})^{-1}b^{(\ell)})_{T+1-\lfloor \ell/2 \rfloor - t} & T - \ell < t \leq T - \lfloor \ell/2 \rfloor \\ 1 & T - \lfloor \ell/2 \rfloor < t \end{cases}$$

$$(17)$$

The definition of $A^{(\ell)} \in \mathbb{R}^{(\ell-\lfloor \ell/2 \rfloor) \times (\ell-\lfloor \ell/2 \rfloor)}$ and $b^{(\ell)} \in \mathbb{R}^{\ell-\lfloor \ell/2 \rfloor}$ is provided in Eq. (22) and Eq. (23) in Appendix A.3.

Moreover, if (X_i, u_i) can only take G values, each with positive probability, $\omega_{g,t} = \frac{1}{|\mathcal{O}_g|} \sum_{i \in \mathcal{O}_g} z_{it}$ is defined similarly as the $\omega_{g,t}$ in Theorem 3. Then any treatment design Z is optimal if it satisfies the constraints in the integer program (4) and

$$\omega_{g,t} = \omega_t^*, \text{ for all } t \text{ and } g.$$
 (18)

Remark 7. Similar as Theorem 3, the dimensions of X_i and u_i (i.e., r and k) can be zero.

Remark 8. To provide a few concrete examples for ω_t^* , if $\ell=1$, we have

$$\omega_t^* = -1 + \frac{2(t-1)}{T-1}$$
 for all t .

If $\ell = 2$, we have

$$\omega_1^* = -1, \quad \omega_2^* = -1 + \frac{2}{2T - 5}, \quad \omega_t^* = -1 + \frac{2t - 3}{T - 2} \quad \text{for } t = 3, \cdots, T - 2, \quad \omega_{T - 1}^* = 1 - \frac{2}{2T - 5}, \quad \omega_T^* = 1.$$

We provide the expression of ω_t^* for $\ell = 3$ in Remark 16 in Appendix A.3. Figure 2 shows the optimal treated proportion for other ℓ in a T = 10 period problem.

REMARK 9. The assumption $T > \frac{\ell^3 + 13\ell^2 + 7\ell + 3}{8\ell}$ can be potentially relaxed. We verified using optimization software that Eq. (17) is optimal even this assumption is not satisfied.

The treated proportion is convex at the beginning and concave towards the end. If the treatment has carryover effects for longer periods (i.e., larger ℓ), the optimal treated proportion is smaller at the beginning and is larger towards the end as shown in Figure 2. There are five different stages in the optimal solution from Theorem 5: in the first stage, all units are control units; in

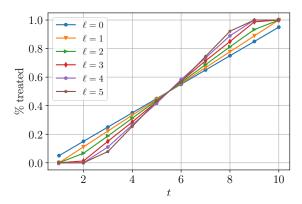


Figure 2 T-optimal design: Optimal treated percentage at each period for a T-period treatment design in the presence of carryover treatment effects, where T = 10. Different colors represent the number of periods ℓ to which the treatment can carry over.

the second stage, ω_t^* grows non-linearly in time and is less than $-1 + \frac{\ell+1}{T-\ell} \frac{t-\lfloor \ell/2 \rfloor}{\lfloor \ell/2 \rfloor+1}$ (implied by the assumption $T > \frac{\ell^3 + 13\ell^2 + 7\ell + 3}{8\ell}$); in the third stage, ω_t^* grows linearly in time; in the fourth stage, ω_t^* grows non-linearly in time again and is larger than $1 - \frac{\ell+1}{T-\ell} \frac{T+\lfloor \ell/2 \rfloor - t}{\lfloor \ell/2 \rfloor + 1}$ (implied by the assumption $T > \frac{\ell^3 + 13\ell^2 + 7\ell + 3}{8\ell}$); in the fifth stage, all units are treated units. The optimal solution is symmetric with respect to the origin, because $f(\omega_1, \dots, \omega_T) = f(-\omega_T, \dots, -\omega_1)$, where $f(\omega_1, \dots, \omega_T)$ is the objective function in the integer program (15).

Intuitively, treated proportions are larger at the end compared with the direct effect only case in Theorem 3 because we want to have sufficiently many observations to precisely estimate the lagged treatment effects. Treated proportions are smaller at the beginning because we only use the observed outcomes from $\ell + 1$ to T to estimate $\vec{\tau}$, α , and $\tilde{\beta}$, so increasing the treated proportion at the beginning can only marginally affect the efficiency.

3.3.1. Overview of the Proof of Theorem 5

The complete proof is deferred to Appendix E and here we provide an overview of the main steps. First, we start with the case where the potential outcome does not have covariates (i.e., r = k = 0 and $\Sigma_e = I_{N(T-\ell)}$) and then maximizing the objective function in the integer program (15) is equivalent to minimizing $\operatorname{tr}(\mathcal{Z}^{\top}\Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top}\mathcal{Z})$. The corresponding convex relaxation is a sum of $\ell+1$ quadratic functions

$$\sum_{j=1}^{\ell+1} \left[\vec{\omega}_{j:(T-\ell-1+j)}^{\top} P_{\vec{1}} \vec{\omega}_{j:(T-\ell-1+j)} + 2(\vec{d}^{(j)})^{\top} \vec{\omega}_{j:(T-\ell-1+j)} \right], \tag{19}$$

where $\vec{\omega}_{j:(T-\ell-1+j)} = \left[\omega_j \cdots \omega_{T-\ell-1+j}\right]^{\top}$, $P_{\vec{1}} = I - \frac{1}{T-\ell} \vec{1} \vec{1}^{\top} \in \mathbb{R}^{(T-\ell)\times(T-\ell)}$, $\vec{d}^{(j)} = [d_t^{(j)}] \in \mathbb{R}^{T-\ell}$, $d_t^{(j)} = \frac{T-\ell-1+2j-2t}{T-\ell}$ for all j and t, and the definition and constraints of ω_t are the same as those of ω_t in the quadratic program (9). The details are presented Lemma 4 in Appendix A.3. Second, we verify

that ω^* in Eq. (17) satisfies the KKT conditions of the quadratic program with objective function (19) and ω^* is therefore an optimal solution for this quadratic program. The details are presented in Theorem 5 in Appendix E. Third, we follow the same steps as Theorem 3 to show additional covariate conditions are needed when the potential outcome has covariates.

4. Algorithms

In this section, we focus on the scenario where the potential outcome has latent covariates, which is widely used in practice,

$$Y(Z) = \alpha \vec{1}^{\top} + \vec{1}\beta^{\top} + UV^{\top} + \tau Z + \varepsilon. \tag{20}$$

For notation simplicity, we use L to denote UV^{\top} in this section.

We propose algorithms to tackle two challenges:

- 1. How can we find the optimal treatment design in practice when covariates are unobserved?
- 2. Is there an alternative (potentially non-linear) estimation approach for τ that may have a smaller variance?

For the first challenge, in Section 4.1, we introduce a local search heuristic (motivated by the sufficient conditions in Theorem 3) to search for a better treatment design by leveraging historical control data that contains latent covariates' information, given that historical control data is commonly available in practice.

For the second challenge, in Section 4.2, we leverage advances on low-rank matrix estimation in the machine learning literature to introduce a non-linear estimator for τ that may be biased but can have a smaller variance. We refer to this estimator as the low-rank matrix estimation with fixed effects (LRME).

4.1. A Data-driven Local Search Heuristic to Find a Better Treatment Design

Since the optimal treatment design depends on the values of latent covariates from Theorem 3, the historical control data with covariates' information could help us find a better treatment design. In order to find a better treatment design, we propose to (a) approximate $Var(\hat{\tau})$ using a robust surrogate when our estimation approach does not a closed-form expression for $Var(\hat{\tau})$, and (b) find an approximate solution to the computationally challenging integer program of finding optimal Z that minimizes $Var(\hat{\tau})$. We propose Algorithm 1 that addresses (a) and (b) as follows:

- For (a), we adopts a robust measure to approximate $Var(\hat{\tau})$ that consists of two phases
- (i) Split the historical control data into sub-matrices: we use a moving window approach to split the historical control data $Y^{\text{hist,All}}$ into m sub-matrices $Y^{\text{hist,(1)}}, \dots, Y^{\text{hist,(m)}}$, each with dimension $N \times T$.

- (ii) We take a synthetic treatment effect τ (can be 0) and apply synthetic experiment to $Y^{\mathrm{hist},(j)}$ using treatment design Z. In more detail, we assume $Y^{\mathrm{hist},(j)}$ does not contain the intervention we study. We define the control outcome as $Y^{(j)}_{it}(-1) = Y^{\mathrm{hist},(j)}_{it}$, the treated outcome as $Y^{(j)}_{it}(1) = Y^{(j)}_{it}(-1) + \tau$, and the "observed" outcome as $Y^{(j)}_{it} = \frac{1+Z_{it}}{2} \cdot Y^{(j)}_{it}(1) + \frac{1-Z_{it}}{2} \cdot Y^{(j)}_{it}(-1)$ for every i, t and j. We use $Y^{(j)}$ and Z as the input to estimate τ and the estimator is denoted as $\hat{\tau}^{(j)}$. We calculate the maximum estimation error across all sub-matrices, $\max_j |\hat{\tau}^{(j)} \tau|$, and use it as a robust surrogate for $\mathrm{Var}(\hat{\tau})$.
- For (b), we start from an arbitrary treatment design Z that satisfies $\frac{1}{N} \sum_{i=1}^{N} z_{it} = \frac{2t-1-T}{T}$ and the constraints in the integer program (4). Then we adopt the following local search heuristic to improve Z (motivated by *simulated annealing*) until Z converges or the algorithm runs for a maximum number of iterations:
- (i) Swap two random rows in Z at every iteration and denote the new treatment design as Z^{new} (note that the cross-sectional average remains unchanged at $\frac{2t-1-T}{T}$ at time period t).
- (ii) Calculate the maximum estimation error $\max_j |\hat{\tau}^{(j)} \tau|$ when treatment design Z^{new} is applied to $Y^{(1)}, \dots, Y^{(m)}$. If $\max_j |\hat{\tau}^{(j)} \tau|$ is reduced, replace Z by Z^{new} . Otherwise, keep Z unchanged with probability $1 p_{\text{escape}}$ but still replace Z by Z^{new} with probability p_{escape} . Here, p_{escape} is a small positive number.
- (iii) At the end of each iteration, define Z^{opt} to be the matrix Z that has the smallest $\max_j |\hat{\tau}^{(j)} \tau|$. Once the algorithm stops (maximum number of iterations is reached or Z is converged), return Z^{opt} .
- Remark 10. We used the maximum estimation error instead of the mean-squared error because it is more robust to outliers.

REMARK 11. In the local search heuristic, when we swap the rows in Z, the cross-sectional average does not change. This guarantees that one of the sufficient conditions for the optimal solution in Theorem 3 (i.e., $\frac{1}{N} \sum_{i=1}^{N} z_{it} = \frac{2t-1-T}{T}$) is always satisfied. Therefore, the goal of the algorithm is to find a treatment design Z that satisfies the other sufficient condition (i.e., $\frac{1}{N} \sum_{i=1}^{N} u_i z_{it} = \mu$).

REMARK 12. In the local search heuristic, we replace Z by Z^{new} with some small probability p_{escape} following the idea of simulated annealing, to escape potential local minima that traps Z.

4.2. An Alternative Estimation Approach for τ

Instead of incorporating L into the error term, we could estimate L jointly with τ using the following objective function,

$$\hat{\tau}, \hat{\alpha}, \hat{\beta}, \hat{L} = \operatorname*{arg\,min}_{\tau, \alpha, \beta, L} \frac{1}{NT} \left\| Y - \alpha \vec{1}^{\top} - \vec{1}\beta^{\top} - L - \tau Z \right\|_{F}^{2} + \mu \left\| L \right\|_{*}, \tag{21}$$

denoted as low-rank matrix estimation with fixed effects (LRME). Here, $||L||_*$ is the nuclear norm (or trace norm) of matrix L, which is equal to the sum of its singular values. Also, $||\cdot||_F$ refers to

Algorithm 1: A data-driven local search heuristic for Z

```
Inputs: Y_{it}^{(1)}(-1), \dots, Y_{it}^{(m)}(-1), Z, \tau, t_{\text{init}}, t_{\text{min}}, v, \text{steps}_{\text{max}}, r, \Delta_{\tau}, t_{\text{max}}
Y_{it}^{(j)}(1) \leftarrow Y_{it}^{(j)}(-1) + \tau \text{ and } Y_{it}^{(j)} \leftarrow \frac{1 + Z_{it}^{\text{current}}}{2} \cdot Y_{it}^{(j)}(1) + \frac{1 - Z_{it}^{\text{current}}}{2} \cdot Y_{it}^{(j)}(-1) \text{ for all } i, \ t \text{ and } j \ ;
for j = 1, \dots, m do
  \hat{\tau}^{(j)}, \hat{\alpha}^{(j)}, \hat{\beta}^{(j)} estimated from Y^{(j)} by OLS/GLS/LRME;
end
\mathbf{E}^{\text{current}} \leftarrow \max_{i} |\hat{\tau}^{(j)} - \tau|, \; \mathbf{E}^{\text{opt}} \leftarrow \mathbf{E}^{\text{current}}, Z^{\text{opt}} \leftarrow Z^{\text{current}}, \; t \leftarrow t_{\text{init}}, \; i_{\text{step}} \leftarrow 1;
while t > t_{\min} and i_{\text{step}} < \text{steps}_{\max} do
         Randomly select distinct i_1, i_2 \in [N];
        \begin{split} Z^{\text{new}} &\leftarrow Z^{\text{current}}, \ Z^{\text{new}}_{i_1,:} \leftarrow Z^{\text{current}}_{i_2,:}, \ Z^{\text{new}}_{i_2,:} \leftarrow Z^{\text{current}}_{i_1,:}, \ i_{\text{step}} \leftarrow i_{\text{step}} + 1 \ ; \\ Y^{(j)}_{it} &\leftarrow \frac{1 + Z^{\text{new}}_{it}}{2} \cdot Y^{(j)}_{it}(1) + \frac{1 - Z^{\text{new}}_{it}}{2} \cdot Y^{(j)}_{it}(-1) \ \text{for all} \ i, \ t \ \text{and} \ j \ ; \end{split}
         for j = 1, \dots, m do
           \hat{\tau}^{(j)}, \hat{\alpha}^{(j)}, \hat{\beta}^{(j)} estimated from Y^{(j)} by OLS/GLS/LRME;
         end
        E^{\text{new}} \leftarrow \max_{i} |\hat{\tau}^{(i)} - \tau|, dE = E^{\text{new}} - E^{\text{current}};
         if E^{new} < E^{opt} then
                 E^{\text{current}} \leftarrow E^{\text{new}}, Z^{\text{current}} \leftarrow Z^{\text{new}}, Z^{\text{opt}} \leftarrow Z^{\text{new}};
                  if \exp(dE/t) > \operatorname{random}(0,1) then E^{\operatorname{current}} \leftarrow E^{\operatorname{new}}, Z^{\operatorname{current}} \leftarrow Z^{\operatorname{new}};
         end
end
Outputs: Z^{\text{opt}}
```

the Frobenius norm of a matrix. If the regularization parameter μ is larger, the rank of \hat{L} tends to be smaller. It is known in the matrix estimation literature that such a bias in estimating L can lead to a lower variance (Candes and Recht 2009, Candes and Plan 2010). If $\text{Var}(\hat{L})$ is lower, $\text{Var}(\hat{\tau})$ is likely to be smaller. We compare LRME with the BLUE estimator through synthetic experiments in Section 5.

The objective function (21) is convex in τ, α, β and L, which has NT + N + T + 1 variables in total. Finding the global optimal solution of convex program (21) can be slow with off-the-shelf software for convex optimization problems such as cvxpy. Alternatively, we propose to use the iterative singular value thresholding and ordinary least squares (iterative SVT and OLS) algorithm to efficiently solve convex program (21). The details of this algorithm are described in Algorithm

2. We can justify using SVT Theorem 1 in Hastie et al. (2015) that shows the optimal solution of

$$\hat{L} = \underset{rank(\ell) < k_0}{\arg \min} \frac{1}{2} \|Y - L\|_F^2 + \mu \|L\|_*,$$

is $\hat{L} = U_{k_0} S_{\mu}(D_{k_0}) V_{k_0}^{\top}$, where the rank- k_0 SVD of Y is $U_{k_0} D_{k_0} V_{k_0}^{\top}$ and $S_{\mu}(D_{k_0})$ is a diagonal k_0 by k_0 matrix with its diagonal entries to be $(\sigma_1 - \mu)_+, \dots, (\sigma_{k_0} - \mu)_+$. When we have historical control data, we can use cross-validation to find the optimal μ by the grid search algorithm. The details are provided in Appendix B.1.

```
Algorithm 2: Iterative SVT and OLS
```

```
\begin{split} &\widehat{\boldsymbol{\tau}}^{(-1)} \leftarrow 0 \;; \\ & \text{At } t = 0, \; \widehat{\boldsymbol{\tau}}^{(0)}, \hat{\boldsymbol{\alpha}}^{(0)}, \hat{\boldsymbol{\beta}}^{(0)} \leftarrow \arg\min_{\tau,\alpha,\beta} \frac{1}{2} \left\| \boldsymbol{Y} - \alpha \vec{\mathbf{1}}^{\top} - \vec{\mathbf{1}} \boldsymbol{\beta}^{\top} - \tau \boldsymbol{Z} \right\|_{F}^{2}; \\ & \hat{\boldsymbol{Y}}_{e}^{(0)} \leftarrow \boldsymbol{Y} - \hat{\boldsymbol{\alpha}}^{(0)} \vec{\mathbf{1}}^{\top} - \vec{\mathbf{1}} (\hat{\boldsymbol{\beta}}^{(0)})^{\top} - \widehat{\boldsymbol{\tau}}^{(0)} \boldsymbol{Z} \;; \\ & \textbf{while} \quad |\widehat{\boldsymbol{\tau}}^{(t)} - \widehat{\boldsymbol{\tau}}^{(t-1)}| > \Delta_{\tau} \; \text{and} \; t < t_{\max} \; \textbf{do} \\ & \quad | \quad \text{The rank-} k_0 \; \text{SVD of} \; \hat{\boldsymbol{Y}}_{e}^{(t)} \; \text{is} \; \boldsymbol{U}_{k_0}^{(t)} \boldsymbol{D}_{k_0}^{(t)} (\boldsymbol{V}_{k_0}^{(t)})^{\top}, \; \text{where} \; \boldsymbol{D}_{k_0}^{(t)} = \operatorname{diag}(\boldsymbol{d}_1^{(t)}, \cdots, \boldsymbol{d}_{k_0}^{(t)}) \;; \\ & \quad \boldsymbol{S}_{\mu_{NT}}(\boldsymbol{D}_{k_0}^{(t)}) \leftarrow \operatorname{diag}((\boldsymbol{d}_1^{(t)} - \mu_{NT})_+, \cdots, (\boldsymbol{d}_{k_0}^{(t)} - \mu_{NT})_+) \;; \\ & \quad \hat{\boldsymbol{L}}^{(t+1)} = \boldsymbol{U}_{k_0}^{(t)} \boldsymbol{S}_{\mu_{NT}}(\boldsymbol{D}_{k_0}^{(t)}) (\boldsymbol{V}_{k_0}^{(t)})^{\top} \;; \\ & \quad \hat{\boldsymbol{\tau}}^{(t+1)}, \hat{\boldsymbol{\alpha}}^{(t+1)}, \hat{\boldsymbol{\beta}}^{(t+1)} = \arg\min_{\tau,\alpha,\beta} \frac{1}{2} \left\| \boldsymbol{Y} - \alpha \vec{\mathbf{1}}^{\top} - \vec{\mathbf{1}} \boldsymbol{\beta}^{\top} - \tau \boldsymbol{Z} - \hat{\boldsymbol{L}}^{(t+1)} \right\|_{F}^{2} \;; \\ & \quad \hat{\boldsymbol{Y}}_{e}^{(t+1)} = \boldsymbol{Y} - \hat{\boldsymbol{\alpha}}^{(t+1)} \vec{\mathbf{1}}^{\top} - \vec{\mathbf{1}} (\hat{\boldsymbol{\beta}}^{(t+1)})^{\top} - \hat{\boldsymbol{\tau}}^{(t+1)} \boldsymbol{Z} \;; \\ & \quad t \leftarrow t + 1 \;; \\ & \quad \textbf{end} \\ & \quad \textbf{Outputs:} \; \hat{\boldsymbol{\tau}}^{(t-1)}, \hat{\boldsymbol{\alpha}}^{(t-1)}, \hat{\boldsymbol{\beta}}^{(t-1)}, \hat{\boldsymbol{L}}^{(t-1)} \end{split}
```

5. Empirical Study: Reducing Influenza Occurrence Rate

Influenza is primarily a community-based infection that is transmitted in households and community settings. CDC estimates that influenza has resulted in 9 million to 45 million illnesses, 140,000 to 810,000 hospitalizations, and 12,000 to 61,000 deaths annually since 2010. The best way to prevent infectious diseases such as seasonal flu is vaccination. However, in the absence of a vaccine, like in the COVID-19 pandemic, other interventions such as social distancing or wearing a face cover are recommended. An important and challenging question is: what is the efficacy of such preventive measures or other interventions in different environments such as schools, workplaces, and businesses in slowing down the spread of the disease and reducing the substantial burden of the disease on society?

The best way to answer this question is to run pilot programs or experiments. Since flu is contagious, we must randomize the experiments at an aggregate geographical level, such as the

Metropolitan Statistical Area (MSA), to minimize the spillover effects. The tradeoff is, we would have a smaller number of units and therefore lower the statistical power to detect the treatment effect. Moreover, the treatment effect is, in general, small compared with the natural fluctuations of flu occurrence rate, so the statistical power is critical to detect the effect of new interventions. Hence, our solution in Section 3 and our algorithm in Section 4 are especially relevant in this setting to evaluate the effect of new interventions.

5.1. Data Description: MarketScan Research Databases

In this study, we focus on the MarketScan Research Databases (formerly called Truven) that contain longitudinal, patient-level, medical claim data for healthcare research. The databases contain data on continuous periods of enrollment per individual, inpatient, and outpatient claim records from the beginning of 2007 to mid-2017. We focus on the inpatient admission records and outpatient service records where the primary diagnosis (DX1) is influenza according to ICD-9-CM diagnosis codes (the databases only have claim records with ICD-9 diagnosis codes). Overall, the data contains 21,277 inpatient admissions and 9,678,572 outpatient records with primary diagnosis influenza, indicating patients are not usually admitted to hospitals for influenza. We consider this combined data for the analysis. We aggregate the influenza records by MSA and month. Furthermore, we aggregate individual enrollment data by MSA and month to get the total number of enrolled patients, by MSA and month. We divide the number of influenza visits by the number of enrolled patients to obtain the flu visit occurrence rate by MSA and month.

Due to the seasonal nature of the flu, we focus on the months from October to April, where flu is active, and there is a relatively high flu visit occurrence rate. We further focus on the period between October 2007 to April 2015 because outside this period, we have only a few observations, and the flu occurrence rate is 0 for most MSAs. We get a panel of 185 MSAs over 56 months, where all MSAs have a positive occurrence rate every month.

REMARK 13. In addition to the flu occurrence problem, we also apply our treatment design algorithm to two other data sets (home medical visits and purchases from a large grocery store). The secondary and tertiary analysis are presented in Appendix F.

5.2. Experiment Setups

In this section, we present the models and estimation methods, treatment designs, how to run synthetic experiments, the metrics to evaluate treatment designs, and the specification used in the synthetic experiments.

 $^{^{5}}$ In ICD-9-CM, the following diagnosis codes indicate influenza: 488, 487.0, 487.1, 487.8, 488.0, 488.1, 488.01, 488.02, 488.09, 488.11, 488.12, 488.19, 488.81, 488.82, and 488.89.

Model Assumptions and Estimation Methods. Since we do not know the data generating process of the flu occurrence data, we consider two models when the treatment only has direct effect: one is the two-way fixed effect model with direct effect only $Y_{it}(z_{it}) = \alpha_i + \beta_t + \tau z_{it} + \varepsilon_{it}$, and the other one is the latent factor model with direct effect only $Y_{it}(z_{it}) = \alpha_i + \beta_t + u_i^{\top} v_t + \tau z_{it} + \varepsilon_{it}$. Moreover, we consider two models with carryover treatment effects: two-way fixed effect model with carryover effects $Y_{it}(z_{it}) = \alpha_i + \beta_t + \tau_1 z_{it} + \cdots + \tau_{\ell+1} z_{i,t-\ell} + \varepsilon_{it}$ and latent factor model with carryover effects $Y_{it}(z_{it}) = \alpha_i + \beta_t + u_i^{\top} v_t + \tau_1 z_{it} + \cdots + \tau_{\ell+1} z_{i,t-\ell} + \varepsilon_{it}$. For the two-way fixed effect model, we use **OLS** to estimate τ (or $\tau_1, \dots, \tau_{\ell+1}$ for the model with carryover effects). For the latent factor model, we use two methods to estimate τ (or $\tau_1, \dots, \tau_{\ell+1}$ for the model with carryover effects): one is feasible **GLS**, 6 and the other one is **LRME** presented in Section 4.2.

Treatment Designs. We compare the different treatment designs via synthetic experiments. Denote ω_t as the cross-sectional average of the treatment design Z at time t and note that,

$$\omega_t = \begin{cases} -1 & \text{all control} \\ 0 & \text{half treated} \\ 1 & \text{all treated} \end{cases}$$

We consider the following treatment designs

- 1. Benchmark treatment designs:
- (a) $Z_{\rm FF}$ (fifty-fifty): $Z_{\rm FF}$ satisfies $\omega_t = 0$ for all t. $Z_{\rm FF}$ is optimal when the potential outcome only has time fixed effects, as shown in Lemma 2.⁷
- (b) Z_{BA} (before-after): Z_{BA} satisfies $\omega_t = -1$ for $t < \frac{T+1}{2}$ and $\omega_t = 1$ for $t > \frac{T+1}{2}$. Z_{BA} is optimal when the potential outcome only has unit fixed effects, as shown in Lemma 3.8
- (c) $Z_{\text{FF+BA}}$ (fifty-fifty with before-after): $Z_{\text{FF+BA}}$ satisfies $\omega_t = -1$ for $t < \frac{T+1}{2}$ and $\omega_t = 0$ for $t \ge \frac{T+1}{2}$, implying half of the units switch from control to treated at midway through the experiment.
- 2. Z_{OPT} : Z_{OPT} satisfies $\omega_t = \frac{2t-1-T}{T}$ for all t. Z_{OPT} is optimal when the potential outcome has both unit and time fixed effects and there are no covariates, as shown in Theorem 3.
- 3. $Z_{K=2}$ (or $Z_{K=3}$, partition and stratification): the treatment design that satisfies $\omega_t = \frac{2t-1-T}{T}$ for each group, where units are partitioned into groups based on their value in the largest singular vectors estimated from the historical data.

⁶ GLS is the BLUE estimator, but is infeasible. We provide the details to implement feasible GLS in Appendix B.2.

⁷ When Z_{FF} satisfies $\omega_t = 0$ for all t, \vec{z} will be in the span of Γ in model (5) so τ can not be uniquely identified. Hence, we randomly select one more unit as control in the first time period and one more unit as treated in the last time period to uniquely identify τ .

⁸ When $Z_{\rm BA}$ satisfies $\omega_t = -1$ for $t < \frac{T+1}{2}$ and $\omega_t = 1$ for $t > \frac{T+1}{2}$, \vec{z} will be in the span of Γ defined in model (5) so τ can not be uniquely identified. Hence, we randomly select one more unit as treated in time period $\frac{T}{2}$ and one more unit as control in time period $\frac{T}{2} + 1$ to uniquely identify τ .

- 4. $Z_{\text{OPT+}}$: the treatment design returned from Algorithm 1, where the input of Algorithm 1 is Z_{OPT} .
- 5. $Z_{\text{OPT-CO}}$: $Z_{\text{OPT-CO}}$ satisfies $\omega_t = \omega_t^*$ for all t, where ω_t^* is defined in Eq. (17) and is nonlinear in t. $Z_{\text{OPT-CO}}$ is optimal when the treatment has carryover effects over ℓ lagged periods and there are no covariates, as shown in Theorem 5.

Synthetic Treatment Effect. The original data does not contain any specific treatment that we can study. We first consider the model with only direct treatment effect. Therefore, entry (i,t)of the original data is $Y_{it}(-1)$. Then, we consider a hypothetical intervention that only affects the current period with synthetic treatment effect τ , and if entry (i,t) is hypothetically assigned to treatment group, i.e., $z_{it} = 1$, we define $Y_{it}(1) = Y_{it}(-1) + \tau$. Therefore, the final input to different estimation methods is going to be the values of matrix Y(Z) as well as Z itself. The value of τ and counterfactual values of each Y(Z) are hidden from the estimators. The true value of τ will later be used to assess performance (RMSE) of the various treatment designs and estimation methods. Moreover, we consider a hypothetical intervention with synthetic carryover treatment effects $\tau_1, \tau_2, \cdots, \tau_{\ell+1}$ that are defined in model (14). We use $Y_{it}((z_{i,t-\ell}, \cdots, z_{i,t}))$ to denote the potential outcome when the past $\ell+1$ periods' treatment status is $z_{i,t-\ell},\cdots,z_{i,t}$. Entry (i,t) of the original data is $Y_{it}((-1,\dots,-1,-1))$. If unit i switches from the control to the treated group at period t, i.e., $z_{i,t-1} = -1, z_{it} = 1$, we define $Y_{it}((-1, \dots, -1, 1)) = Y_{it}((-1, \dots, -1, -1)) + \tau_1$; If unit i switches from the control to the treated group at period t-s $(1 \le s < \ell)$, we define $Y_{it}((-1,\dots,1,\dots,1)) = Y_{it}((-1,\dots,-1,\dots,-1)) + \tau_1 + \tau_2 + \dots + \tau_{s+1};$ If unit *i* switches from the control to the treated group at period t-s $(s \ge \ell)$, we define $Y_{it}((1, \dots, 1)) = Y_{it}((-1, \dots, -1)) + (-1, \dots, -1)$ $\tau_1+\tau_2+\cdots+\tau_{\ell+1}.$

Evaluation Metrics. We randomly select m blocks each with dimension $N \times \tilde{T}$. For each block, the first $\tilde{T} - T$ time periods are the historical control data and we apply synthetic experiments to the last $N \times T$ time periods. We split the historical control data into k matrices, where two adjacent matrices could have overlapping time periods. We use these k matrices as the input for Algorithm 1. In the experiment with direct treatment effect only, for each treatment design Z and each of the m blocks, the estimated treatment effect is denoted as $\hat{\tau}_Z^{(j)}$, where $j=1,\cdots,m$. We report the following metrics for each treatment design Z:

$$\text{RMSE}_{Z} = \left(\frac{1}{m} \sum_{j=1}^{m} \left(\hat{\tau}_{Z}^{(j)} - \tau\right)^{2}\right)^{1/2}, \quad \bar{\tau}_{Z} = \frac{1}{m} \sum_{j=1}^{m} \hat{\tau}_{Z}^{(j)}, \quad \text{Var}(\hat{\tau}_{Z}) = \frac{1}{m} \sum_{j=1}^{m} \left(\hat{\tau}_{Z}^{(j)} - \bar{\tau}_{Z}\right)^{2}.$$

In the experiment with carryover effects, the estimated direct and lagged *i*-period treatment effects are denoted as $\hat{\tau}_{Z,1}^{(j)}$ and $\hat{\tau}_{Z,i}^{(j)}$ for $i=1,\cdots,\ell+1$ and $j=1,\cdots,m$. We report the RMSE across all treatment effect parameters

$$\text{RMSE}_{Z,\text{Carryover}} = \left(\frac{1}{m(\ell+1)} \sum_{j=1}^{m} \sum_{l=1}^{\ell+1} \left(\hat{\tau}_{Z,l}^{(j)} - \tau_l\right)^2\right)^{1/2}.$$

Specification. We choose the number of random blocks at m = 500 and the treatment effect at $\tau = -0.01\%$ in the experiment with direct treatment effect only. Note that the median occurrence rate is 0.0714%, so the treatment effect is about 14% of the median occurrence rate. There are 61 MSAs with at least 0.01%. occurrence rate in every month. Our synthetic experiments are implemented on these 61 MSAs. Moreover, in the experiment with carryover effects, we consider the prevention policy can affect the outcome for the current period and two periods in the future. We choose the treatment effects at $[\tau_1, \tau_2, \tau_3] = [-0.007\%, -0.002\%, -0.001\%]$ (so the occurrence rate stays positive). It is natural to start a new flu prevention policy at the beginning of a flu season, roll out the new policy during the flu reason, and terminate the policy at the end of the flu season (probably restart the policy at the beginning of the next flu season). Hence, we use T = 7, start the experiment in October and end in April. We use N = 25 and N = 50.

5.3. Results

In this section, we demonstrate three main findings from the synthetic experiments on the flu data as shown in Tables 4-5 and Figure 3 and on the home medical visit data and grocery data as shown in Tables 6-9 and Figures 5-6 in Appendix F:

- 1. The linear staggered treatment design $Z_{\rm OPT}$ can significantly outperform all benchmark treatment designs.
- 2. We can use historical data to search for a better treatment design compared with Z_{OPT} , for example, using our data-driven local search heuristic.
- 3. When the treatment has carryover effects, the nonlinear staggered treatment design $Z_{\text{OPT-CO}}$ outperforms the linear treatment design Z_{OPT} .

Linear Staggered Treatment Design Outperforms Benchmark Treatment Designs. We first compare the RMSE of the optimal treatment design matrix Z_{OPT} with the benchmark treatment designs Z_{FF} , Z_{BA} , and $Z_{\text{FF+BA}}$ for all estimation methods in Table 4. Figure 3 shows the bias and variance of benchmark treatment designs and Z_{OPT} on the flu data. Tables 6 and 8 and Figures 5 and 6 in Appendix F show the corresponding results for the home medical visit data and grocery data. We have the following findings:

1. Z_{OPT} can consistently outperform benchmark treatment designs for all estimation methods.

| | (N,T) | (25,7) | (50,7) | | (N,T) | (25,7) | (50,7) |
|----------------|----------------------|---------|---------|----------------|---------------------|---------|---------|
| OLS | Z_{BA} | 0.14273 | 0.14311 | OLS | $Z_{ m OPT}$ | 0.03539 | 0.02530 |
| | $Z_{ m FF}$ | 0.05939 | 0.06630 | | $Z_{K=2}$ | 0.02916 | 0.02005 |
| | $Z_{\text{FF+BA}}$ | 0.03677 | 0.02631 | | $Z_{K=3}$ | 0.02894 | 0.01920 |
| | $Z_{ m OPT}$ | 0.03539 | 0.02530 | | $Z_{\mathrm{OPT+}}$ | 0.02989 | 0.01990 |
| GLS | Z_{BA} | 0.13648 | 0.13891 | GLS | $Z_{ m OPT}$ | 0.02290 | 0.01524 |
| | $Z_{ m FF}$ | 0.02871 | 0.02832 | | $Z_{K=2}$ | 0.02070 | 0.01303 |
| | $Z_{\mathrm{FF+BA}}$ | 0.02771 | 0.01894 | | $Z_{K=3}$ | 0.02004 | 0.01279 |
| | $Z_{ m OPT}$ | 0.02290 | 0.01524 | | $Z_{\mathrm{OPT+}}$ | 0.02104 | 0.01285 |
| LRME | Z_{BA} | 0.12495 | 0.12508 | LRME | $Z_{ m OPT}$ | 0.02188 | 0.01568 |
| | Z_{FF} | 0.02909 | 0.03120 | | $Z_{K=2}$ | 0.02221 | 0.01482 |
| | $Z_{\text{FF+BA}}$ | 0.02323 | 0.01659 | | $Z_{K=3}$ | 0.02184 | 0.01512 |
| | Z_{OPT} | 0.02188 | 0.01568 | | $Z_{\mathrm{OPT+}}$ | 0.01990 | 0.01403 |
| Hist Winner | $Z_{ m OPT}$ | 0.02279 | 0.01558 | Hist Winner | $Z_{ m OPT+}$ | 0.02071 | 0.01451 |

Table 4 Inpatient and outpatient flu visit rates: This table compares the RMSE based on m = 500 randomly sampled blocks for benchmark treatment designs, $Z_{\rm OPT}$, $Z_{K=2}$, $Z_{K=3}$, $Z_{\rm OPT+}$ and "Hist Winner" (the estimation method with minimax estimation error on historical matrices). In the synthetic experiments, the intervention only has direct treatment effect and we choose the treatment effect at $\tau = -0.01\%$. The left table shows the linear staggered design $Z_{\rm OPT}$ outperforms benchmark treatment designs. The right table shows we can find a better treatment design $Z_{\rm OPT+}$ via historical data and our data-driven local search algorithm. The optimal estimation method found on historical data has similar performance as GLS and LRME.

| | (N,T) | (25,7) | (50,7) |
|------|-----------------|--|---------------------------|
| OLS | $Z_{ m OPT-CO}$ | $\begin{vmatrix} 0.04566 \\ 0.04228 \end{vmatrix}$ | $0.03472 \\ 0.03141$ |
| GLS | $Z_{ m OPT}$ | 0.03350 0.03017 | 0.02564 0.02321 |
| LRME | $Z_{ m OPT-CO}$ | $\begin{vmatrix} 0.03170 \\ 0.03104 \end{vmatrix}$ | 0.02508 0.02449 |

Table 5 Inpatient and outpatient $\overline{\text{flu}}$ visit rates: This table compares the RMSE based on m = 500 randomly sampled blocks for Z_{OPT} and $Z_{\text{OPT-CO}}$. In the synthetic experiments, the intervention has direct and carryover treatment effects and we choose the treatment effects at $[\tau_1, \tau_2, \tau_3] = [-0.007\%, -0.002\%, -0.001\%]$. This table shows the nonlinear staggered design $Z_{\text{OPT-CO}}$ outperforms the linear staggered design Z_{OPT} .

2. The latent factor model estimated from either GLS or LRME outperforms the two-way fixed effect model estimated from OLS on all data sets. In other words, GLS and LRME are better estimation methods compared with OLS.

REMARK 14. In practice, we can use cross-validation to select the best estimation method using historical control data. Specifically, we use the leave one (historical matrix) out cross-validation: First, apply synthetic treatment to every historical matrix and estimate τ by OLS, GLS, and LRME; Second, calculate the maximum estimation error of τ across all historical matrices; Third, find the estimation method with the smallest maximum estimation error. We use this method to

estimate the treatment effect on experimental data. We find that the RMSE from this approach is very close to the optimal RMSE, as shown in Table 4.

An Even Better Treatment Design Found via Historical Data. Suppose the potential outcome has either observed or latent covariates. In that case, Theorem 3 shows that the treated proportion of Z_{OPT} satisfies one of the two sufficient conditions for the optimal treatment design. We can apply Algorithm 1 or K-means to the historical data that contains information about covariates to search for a better treatment design that mimics stratification and satisfies the other sufficient condition.

Table 4 presents the RMSE of $Z_{\rm OPT}$, $Z_{K=2}$ and $Z_{K=3}$ returned from K-means, and $Z_{\rm OPT+}$ returned from Algorithm 1 for all estimation methods. Figure 3 shows the bias and variance of these four treatment designs on the flu data. The corresponding results for the other two data sets are presented in Tables 6 and 8 and Figures 5 and 6 in Appendix F. For all estimation methods and all data sets, $Z_{\rm OPT+}$, $Z_{K=2}$, and $Z_{K=3}$ have a smaller RMSE compared with $Z_{\rm OPT}$. Moreover, in general, $Z_{\rm OPT+}$ has a smaller RMSE copmared with $Z_{K=2}$ and $Z_{K=3}$, especially on the medical home visit and grocery data. In empirical applications, we may not have the information about how many groups to partition units into, i.e., K=2,3, or some larger number. Thus, we suggest using $Z_{\rm OPT+}$.

Nonlinear Staggered Treatment Design Outperforms When the Treatment Has Carryover Effects. When the treatment can affect multiple periods, Theorem 5 shows that the optimal treated proportion is nonlinear in time and equals the treated proportion of $Z_{\rm OPT-CO}$. $Z_{\rm OPT-CO}$ can consistently outperform $Z_{\rm OPT}$ for all estimation methods, as shown in Table 5. Moreover, GLS is a better estimation method compared with OLS and LRME. The findings are robust to the other two data sets, as shown in Tables 7 and 9.

6. Concluding Remarks

In this paper, we study the optimal multi-period design of experiments, where our goal is to maximize the statistical power of our estimate for the treatment effect. Formally, we minimize the variance of the estimated treatment effect. We assume that the treated units cannot switch back to control during the experiment. In the simplest case, we assume the potential outcome follows a two-way fixed-effect linear model; In extensions, we consider additional observed or latent covariates and the treatment with carryover effects. We find a set of sufficient conditions for the optimal treatment design. When the potential outcome does not have covariates or has covariates that take discrete values, we can provide a feasible analytical solution that approximates the global optimum within a factor of $1 + O\left(\frac{1}{N^2}\right)$ based on the sufficient conditions.

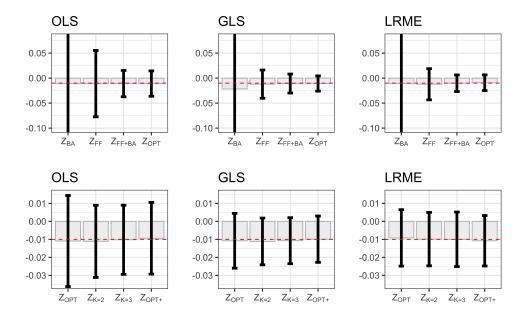


Figure 3 Inpatient and outpatient flu visit rates: This figure compares the average estimated treatment effect $\hat{\tau}_Z$ and its standard deviation $\sqrt{\mathrm{Var}(\hat{\tau}_Z)}$ based on m=500 randomly sampled blocks for benchmark designs, $Z_{\mathrm{OPT}}, Z_{K=2}, Z_{K=3}$ and $Z_{\mathrm{OPT+}}$. In the synthetic experiments, N=50, T=7, the intervention only has direct treatment effect and we choose the treatment effect at $\tau=-0.01\%$. The height of the bar shows $\hat{\tau}_Z$, while the error bar indicates the standard deviation $\sqrt{\mathrm{Var}(\hat{\tau}_Z)}$. The red dash line indicates the true value of $\tau=-0.01\%$. Note that figures in the second row have a different y-axis scale due to superior performance of $Z_{\mathrm{OPT}}, Z_{K=2}, Z_{K=3}$ and $Z_{\mathrm{OPT+}}$ over benchmark treatment designs. The bias of various treatment designs is similar while Z_{OPT} has much smaller variance compared with benchmark designs and the treatment designs found using historical data, such as $Z_{\mathrm{OPT+}}$ from Algorithm 1, has smaller variance compared with $Z_{\mathrm{OPT-}}$.

If we have historical control data, we propose a local search algorithm based on the minimax decision rule to find a better treatment design based on the sufficient conditions when the potential outcome has latent covariates or covariates can take infinitely many values. In the end, we demonstrate the practical relevance of our results through synthetic interventions on the influenza occurrence rate and synthetic experiments on in-home medical service data and grocery data. Our solution can significantly outperform benchmark treatment designs independent of the model assumption and treatment effect estimation approach.

Managerial Implications. Our results have several implications on guiding companies and policymakers to conduct experiments for evaluating interventions. First, when the sample size is small, by leveraging the time dimension and tracking the same sample across multiple time periods, one can increase the number of observations and use our optimal multi-period experiment design to maximize the benefit of the experiment. Second, structure of the multi-period experiment design is substantially impacted when the effect of interventions last for multiple time periods. Third, in

presence of historical (pre-experimentation) data, we can further optimize our treatment designs and hence reduce the number of required samples which means lowering the experiment cost.

Future directions. One interesting direction is how we can dynamically adjust the treatment design to improve the efficiency using both experimental and historical control data, under the constraint that the treatment cannot be removed once implemented, or it is costly to switch from treated to control. This direction is related to online learning or multi-arm/contextual bandit problems, but the objective is different, and there is a switching constraint from treated to control.

References

- Abadie, Alberto, Alexis Diamond, Jens Hainmueller. 2010. Synthetic control methods for comparative case studies: Estimating the effect of californias tobacco control program. *Journal of the American statistical Association* **105**(490) 493–505.
- Aral, Sinan, Dylan Walker. 2014. Tie strength, embeddedness, and social influence: A large-scale networked experiment. *Management Science* **60**(6) 1352–1370.
- Assmus, Gert, John U Farley, Donald R Lehmann. 1984. How advertising affects sales: Meta-analysis of econometric results. *Journal of Marketing Research* **21**(1) 65–74.
- Athey, S, S Stern. 2002. The impact of information technology on emergency health care outcomes. *The Rand Journal of Economics* **33**(3) 399.
- Athey, Susan, Mohsen Bayati, Nikolay Doudchenko, Guido Imbens, Khashayar Khosravi. 2018. Matrix completion methods for causal panel data models. Tech. rep., National Bureau of Economic Research.
- Athey, Susan, Guido W Imbens. 2017. The econometrics of randomized experiments. *Handbook of Economic Field Experiments*, vol. 1. Elsevier, 73–140.
- Athey, Susan, Guido W Imbens. 2018. Design-based analysis in difference-in-differences settings with staggered adoption. Tech. rep., National Bureau of Economic Research.
- Athey, Susan, Scott Stern. 1998. An empirical framework for testing theories about complimentarity in organizational design. Tech. rep., National Bureau of Economic Research.
- Bai, Jushan. 2003. Inferential theory for factor models of large dimensions. Econometrica 71(1) 135–171.
- Bai, Jushan, Serena Ng. 2002. Determining the number of factors in approximate factor models. *Econometrica* **70**(1) 191–221.
- Baltagi, Badi. 2008. Econometric analysis of panel data. John Wiley & Sons.
- Bertsimas, Dimitris, Mac Johnson, Nathan Kallus. 2015. The power of optimization over randomization in designing experiments involving small samples. *Operations Research* **63**(4) 868–876.
- Bertsimas, Dimitris, Nikita Korolko, Alexander M Weinstein. 2019. Covariate-adaptive optimization in online clinical trials. $Operations\ Research$.

- Bhat, Nikhil, Vivek F Farias, Ciamac C Moallemi, Deeksha Sinha. 2019. Near optimal ab testing. *Management Science*.
- Brown, Celia A, Richard J Lilford. 2006. The stepped wedge trial design: a systematic review. BMC medical research methodology $\mathbf{6}(1)$ 54.
- Brynjolfsson, Erik, Avinash Collis, Felix Eggers. 2019. Using massive online choice experiments to measure changes in well-being. *Proceedings of the National Academy of Sciences* **116**(15) 7250–7255.
- Brynjolfsson, Erik, Felix Eggers, Avinash Gannamaneni. 2018. Measuring welfare with massive online choice experiments: A brief introduction. *AEA Papers and Proceedings*, vol. 108, 473–76.
- Bubeck, Sébastien, Nicolo Cesa-Bianchi. 2012. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. arXiv preprint arXiv:1204.5721.
- Candes, Emmanuel, Benjamin Recht. 2009. Exact matrix completion via convex optimization. Communications of the ACM 55(6) 111–119.
- Candes, Emmanuel J, Yaniv Plan. 2010. Matrix completion with noise. *Proceedings of the IEEE* **98**(6) 925–936.
- Cheng, Russell CH, Teresa Davenport. 1989. The problem of dimensionality in stratified sampling. *Management Science* **35**(11) 1278–1296.
- Cheung, Wang Chi, David Simchi-Levi, He Wang. 2017. Dynamic pricing and demand learning with limited price experimentation. *Operations Research* **65**(6) 1722–1731.
- Cohen, Jacob. 1992. Statistical power analysis. Current directions in psychological science 1(3) 98–101.
- Cohen, Maxime, Michael-David Fiszer, Baek Jung Kim. 2018. Frustration-based promotions: Field experiments in ride-sharing. $Available\ at\ SSRN\ 3129717$.
- Cook, Thomas D, Donald Thomas Campbell, William Shadish. 2002. Experimental and quasi-experimental designs for generalized causal inference. Houghton Mifflin Boston, MA.
- Crowder, Martin J, David J Hand. 1990. Analysis of repeated measures, vol. 41. CRC Press.
- Cui, Ruomeng, Jun Li, Dennis J Zhang. 2020. Reducing discrimination with reviews in the sharing economy: Evidence from field experiments on airbnb. *Management Science* **66**(3) 1071–1094.
- Fan, Jianqing, Kosuke Imai, Han Liu, Yang Ning, Xiaolin Yang. 2016. Improving covariate balancing propensity score: A doubly robust and efficient approach. Tech. rep., Technical report, Princeton Univ.
- Fogarty, Colin B, Dylan S Small. 2016. Sensitivity analysis for multiple comparisons in matched observational studies through quadratically constrained linear programming. *Journal of the American Statistical Association* **111**(516) 1820–1830.
- Fox, Bennet L. 2000. Separability in optimal allocation. Operations Research 48(1) 173-176.
- Freeman, PR. 1989. The performance of the two-stage analysis of two-treatment, two-period crossover trials. Statistics in medicine 8(12) 1421–1432.

- Gordon, Brett R, Florian Zettelmeyer, Neha Bhargava, Dan Chapsky. 2019. A comparison of approaches to advertising measurement: Evidence from big field experiments at facebook. *Marketing Science* **38**(2) 193–225.
- Grizzle, James E. 1965. The two-period change-over design and its use in clinical trials. Biometrics 467–480.
- Hastie, Trevor, Rahul Mazumder, Jason D Lee, Reza Zadeh. 2015. Matrix completion and low-rank svd via fast alternating least squares. The Journal of Machine Learning Research 16(1) 3367–3402.
- Hayes, Brian. 2002. Computing science: The easiest hard problem. American Scientist 90(2) 113-117.
- Hemming, Karla, Terry P Haines, Peter J Chilton, Alan J Girling, Richard J Lilford. 2015. The stepped wedge cluster randomised trial: rationale, design, analysis, and reporting. *Bmj* **350** h391.
- Heng, Siyu, Hyunseung Kang, Dylan S Small, Colin B Fogarty. 2019. Increasing power for observational studies of aberrant response: An adaptive approach. $arXiv\ preprint\ arXiv:1907.06770$.
- Hills, M, P Armitage. 1979. The two-period cross-over clinical trial. *British journal of clinical pharmacology* 8(1) 7.
- Hsiao, Cheng. 2014. Analysis of panel data. 54, Cambridge university press.
- Hu, Yanqing, Feifang Hu, et al. 2012. Asymptotic properties of covariate-adaptive randomization. *The Annals of Statistics* **40**(3) 1794–1815.
- Hussey, Michael A, James P Hughes. 2007. Design and analysis of stepped wedge cluster randomized trials. Contemporary Clinical Trials 28(2) 182–191.
- Imai, Kosuke, Marc Ratkovic. 2014. Covariate balancing propensity score. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **76**(1) 243–263.
- Jones, Byron, Michael G Kenward. 2014. Design and analysis of cross-over trials. CRC press.
- Joo, Mingyu, Michael L Thompson, Greg M Allenby. 2019. Optimal product design by sequential experiments in high dimensions. *Management Science* **65**(7) 3235–3254.
- Kallus, Nathan. 2018. Optimal a priori balance in the design of controlled experiments. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 80(1) 85–112.
- Kallus, Nathan. 2020. Generalized optimal matching methods for causal inference. *Journal of Machine Learning Research* **21**(62) 1–54.
- Kallus, Nathan, Xiaojie Mao, Madeleine Udell. 2018. Causal inference with noisy and missing covariates via matrix factorization. Advances in neural information processing systems. 6921–6932.
- Krieger, Abba M, David Azriel, Adam Kapelner. 2019. Nearly random designs with greatly improved balance. Biometrika 106(3) 695–701.
- Kuhfeld, Warren F. 2005. Experimental design, efficiency, coding, and choice designs. *Marketing Research Methods in SAS: Experimental Design, Choice, Conjoint, and Graphical Techniques* 47–97.

- Lattimore, Tor, Csaba Szepesvári. 2018. Bandit algorithms. preprint 28.
- Lawrie, Jock, John B Carlin, Andrew B Forbes. 2015. Optimal stepped wedge designs. *Statistics & Probability Letters* **99** 210–214.
- Leahey, Tricia M, Graham Thomas, Joseph L Fava, Leslee L Subak, Michael Schembri, Katie Krupel, Rajiv Kumar, Brad Weinberg, Rena R Wing. 2014. Adding evidence-based behavioral weight loss strategies to a statewide wellness campaign: a randomized clinical trial. *American Journal of Public Health* 104(7) 1300–1306.
- Lee, Dokyun, Kartik Hosanagar. 2019. How do recommender systems affect sales diversity? a cross-category investigation via randomized field experiment. *Information Systems Research* **30**(1) 239–259.
- Leone, Robert P. 1995. Generalizing what is known about temporal aggregation and advertising carryover.

 Marketing Science 14(3_supplement) G141–G150.
- Li, Fan, Kari Lock Morgan, Alan M Zaslavsky. 2018a. Balancing covariates via propensity score weighting. Journal of the American Statistical Association 113(521) 390–400.
- Li, Fan, Elizabeth L Turner, John S Preisser. 2018b. Optimal allocation of clusters in cohort stepped wedge designs. Statistics & Probability Letters 137 257–263.
- Lu, Yina, Andrés Musalem, Marcelo Olivares, Ariel Schilkrut. 2013. Measuring the effect of queues on customer purchases. *Management Science* **59**(8) 1743–1763.
- lEcuyer, Pierre, Patrick Maillé, Nicolás E Stier-Moses, Bruno Tuffin. 2017. Revenue-maximizing rankings for online platforms with quality-sensitive consumers. *Operations Research* **65**(2) 408–423.
- Mertens, Stephan. 2006. The easiest hard problem: Number partitioning. Computational Complexity and Statistical Physics 125(2) 125–139.
- Muchnik, Lev, Sinan Aral, Sean J Taylor. 2013. Social influence bias: A randomized experiment. *Science* **341**(6146) 647–651.
- Mulvey, John M. 1983. Multivariate stratified sampling by optimization. *Management Science* **29**(6) 715–724.
- Nikolaev, Alexander G, Sheldon H Jacobson, Wendy K Tam Cho, Jason J Sauppe, Edward C Sewell. 2013. Balance optimization subset selection (boss): An alternative approach for causal inference with observational data. *Operations Research* **61**(2) 398–412.
- Panniello, Umberto, Michele Gorgoglione, Shawndra Hill, Kartik Hosanagar. 2014. Incorporating profit margins into recommender systems: A randomized field experiment of purchasing behavior and consumer trust.
- Pocock, Stuart J, Richard Simon. 1975. Sequential treatment assignment with balancing for prognostic factors in the controlled clinical trial. *Biometrics* 103–115.
- Polit, Denise F, Cheryl Tatano Beck. 2008. Nursing research: Generating and assessing evidence for nursing practice. Lippincott Williams & Wilkins.

- Rosenbaum, Paul R, Donald B Rubin. 1983. The central role of the propensity score in observational studies for causal effects. Biometrika **70**(1) 41–55.
- Salganik, Matthew. 2019. Bit by bit: Social research in the digital age. Princeton University Press.
- Sauppe, Jason J, Sheldon H Jacobson. 2017. The role of covariate balance in observational studies. *Naval Research Logistics (NRL)* **64**(4) 323–344.
- Sauppe, Jason J, Sheldon H Jacobson, Edward C Sewell. 2014. Complexity and approximation results for the balance optimization subset selection model for causal inference in observational studies. *INFORMS Journal on Computing* **26**(3) 547–566.
- Senn, Stephen S. 2002. Cross-over trials in clinical research, vol. 5. John Wiley & Sons.
- Simon, Richard. 1979. Restricted randomization designs in clinical trials. *Biometrics* 503–512.
- Sun, Lei, Alexander G Nikolaev. 2016. Mutual information based matching for causal inference with observational data. The Journal of Machine Learning Research 17(1) 6990–7020.
- Sun, Tianshu, Siva Viswanathan, Ni Huang, Elena Zheleva. 2020a. Designing promotional incentive to embrace social sharing: Evidence from field and lab experiments. MIS Quarterly.
- Sun, Tianshu, Siva Viswanathan, Elena Zheleva. 2020b. Creating social contagion through firm-mediated message design: Evidence from a randomized field experiment. *Management Science*.
- Ülkü, Sezer, Chris Hydock, Shiliang Cui. 2019. Making the wait worthwhile: Experiments on the effect of queueing on consumption. *Management Science*.
- Wallenstein, Sylvan, Alan C Fisher. 1977. The analysis of the two-period repeated measurements crossover design with application to clinical trials. *Biometrics* 261–269.
- Willan, Andrew R, Joseph L Pater. 1986. Carryover and the two-period crossover clinical trial. *Biometrics* 593–599.
- Woertman, Willem, Esther de Hoop, Mirjam Moerbeek, Sytse U Zuidema, Debby L Gerritsen, Steven Teerenstra. 2013. Stepped wedge designs could reduce the required sample size in cluster randomized trials.

 *Journal of Clinical Epidemiology 66(7) 752–758.
- Zhang, Dennis J, Hengchen Dai, Lingxiu Dong, Fangfang Qi, Nannan Zhang, Xiaofei Liu, Zhongyi Liu, Jiang Yang. 2019a. The long-term and spillover effects of price promotions on retailing platforms: Evidence from a large randomized experiment on alibaba. *Management Science*.
- Zhang, Dennis J, Hengchen Dai, Lingxiu Dong, Qian Wu, Lifan Guo, Xiaofei Liu. 2019b. The value of popup stores on retailing platforms: Evidence from a field experiment with alibaba. *Management Science* **65**(11) 5142–5151.
- Zhao, Qingyuan. 2019. Covariate balancing propensity score by tailored loss functions. *The Annals of Statistics* **47**(2) 965–993.

Study Approval and Acknowledgement: This paper used the MarketScan Research Databases. Data access was provided by the Stanford Center for Population Health Sciences Data Core. The PHS Data Core is supported by a National Institutes of Health National Center for Advancing Translational Science Clinical and Translational Science Award (UL1 TR001085) and internal Stanford funding. The content is solely the responsibility of the authors and does not necessarily

Appendix A: Additional Theoretical Results

A.1. Treatment is Reversible

represent the official views of the NIH.

PROPOSITION 1. Suppose Assumption 1 holds, and the treatment is reversible, let $\omega_t = \frac{1}{N} \sum_{i=1}^{N} z_{it}$ and $\zeta_i = \frac{1}{T} \sum_{i=1}^{T} z_{it}$.

1. Suppose the potential outcome can be modeled as $Y_{it}(z_{it}) = \beta_t + \tau z_{it} + \varepsilon_{it}$. Any treatment design is optimal if it satisfies

$$\omega_{t} = 0.$$

2. Suppose the potential outcome can be modeled as $Y_{it}(z_{it}) = \alpha_i + \tau z_{it} + \varepsilon_{it}$. Any treatment design is optimal if it satisfies

$$\zeta_i = 0.$$

3. Suppose the potential outcome can be modeled as $Y_{it}(z_{it}) = \alpha_i + \beta_t + \tau z_{it} + \varepsilon_{it}$. Any treatment design is optimal if it satisfies

$$\omega_t = 0, \quad \zeta_i = 0.$$

4. Suppose the potential outcome can be modeled as $Y_{it}(z_{it}) = \alpha_i + \beta_t + X_i^{\top}\theta_t + u_i^{\top}v_t + \tau z_{it} + \varepsilon_{it}$ (or $Y_{it}(z_{it}) = \alpha_i + \beta_t + X_i\theta_t + \tau_1 z_{it} + \tau_2 z_{i,t-1} + \cdots + \tau_{\ell+1} z_{i,t-\ell} + u_i^{\top}v_t + \varepsilon_{it}$). Suppose the assumptions in Theorem 3 (or Theorem 5) hold and let $\omega_{g,t} = \frac{1}{|\mathcal{O}_g|} \sum_{i \in \mathcal{O}_g} z_{it}$, where $\mathcal{O}_g = \{i : X_i = x_g, u_i = u_{0,g}\}$. Any treatment design is optimal if it satisfies

$$\omega_{q,t} = 0$$
, for all t and g, $\zeta_i = 0$.

Proof of Proposition 1 The proof is a direct extension of the proof of Lemmas 2 and 3, Theorems 1, 3 and 5 in Appendix C. The major difference is when the treatment can be reversed, the value of $\sum_{i=1}^{N} \zeta_i^2$ is not unique given ω_t . Then we can separately optimize the part with ζ_i and ω_t . Note that $\sum_{i=1}^{N} \zeta_i^2$ is minimized at $\zeta_i = 0$. Since the Hessian of the part with ω_t ($\omega_{g,t}$) is positive semi-definite, this part is minimized at $\omega_t = 0$ (and $\omega_{g,t} = 0$). \square

A.2. Additional Discussion of Theorem 3

REMARK 15. Consider a symmetric case for Theorem 3 that rows in V are deterministic and have $\sum_{i=1}^{T} v_i = \vec{0}$, u_i is random with mean 0 and variance I_k , $Cov(u_i, u_j) = 0$ for any $i \neq j$. In this case, errors

are cross-sectional independent but time-series correlated, and units are exchangeable. By symmetricity, we have the sufficient conditions for the optimal solution in this case (similar to Theorem 3),

$$\frac{1}{N} \sum_{i=1}^{N} z_{it} = \frac{2t - 1 - T}{T}, \quad \frac{1}{N} \sum_{i=1}^{N} X_i z_{it} = \mu_X, \quad \frac{1}{T} \sum_{t=1}^{T} v_t z_{it} = \mu_V, \text{ for all } i$$

for some $\mu_V \in \mathbb{R}^k$. When z_{it} is the same for all i (before-after design Z_{BA} defined in Section 5.2), the time-series average condition $\frac{1}{T}\sum_{t=1}^T v_t z_{it} = \mu_V$ for all i is satisfied. However, this treatment design violates the first condition on the cross-sectional average, $\frac{1}{N}\sum_{i=1}^N z_{it} = \frac{2t-1-T}{T}$. When v_t is not stationary and has a strong time trend, such as GDP and stock indices over time, there may not exist a Z to satisfy all the sufficient conditions. Intuitively, the optimal Z balances the sufficient conditions on the time-series and cross-sectional average, which is closer to Z_{BA} rather than Z_{FF} .

A.3. More Results for Carryover Treatment Effects

In Theorem 5, $A^{(\ell)}$ and $b^{(\ell)}$ in Eq. (17) are defined as follows

$$A^{(\ell)} = \begin{bmatrix} \lfloor \ell/2 \rfloor + 1 & & \\ & \lfloor \ell/2 \rfloor + 2 & \\ & & \ddots & \\ & & & \ell \end{bmatrix} - \frac{1}{T - \ell} \begin{bmatrix} \ell - \lfloor \ell/2 \rfloor & \ell - 1 - \lfloor \ell/2 \rfloor & \ell - 2 - \lfloor \ell/2 \rfloor & \cdots & 1 \\ \ell - 1 - \lfloor \ell/2 \rfloor & \ell - 1 - \lfloor \ell/2 \rfloor & \ell - 2 - \lfloor \ell/2 \rfloor & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}$$
(22)

$$b^{(\ell)} = -\begin{bmatrix} \lfloor \ell/2 \rfloor + 1 \\ \vdots \\ \ell - 1 \\ \ell \end{bmatrix} + \frac{1}{T - \ell} \begin{bmatrix} (\lfloor \ell/2 \rfloor + 1)^2 \\ \vdots \\ (\ell - 1)^2 \\ \ell^2 \end{bmatrix} - \frac{1}{T - \ell} \begin{bmatrix} \sum_{l=1}^{\ell - \lfloor \ell/2 \rfloor} (\lfloor \ell/2 \rfloor + 1 - l) \\ \vdots \\ 2\lfloor \ell/2 \rfloor - 1 \\ \lfloor \ell/2 \rfloor \end{bmatrix}$$
(23)

The first step to show Theorem 5 is to provide the convex relaxation of the integer program (15) when the potential outcome does not have observed or latent covariates. When the model does not have latent covariates, $\Sigma_e = \sigma^2 I_{N(T-\ell)}$ and $\text{Var}(\hat{\vec{\tau}})$ equals

$$\operatorname{Var}(\hat{\vec{\tau}}) = (\mathcal{Z}^{\top} (I_{N(T-\ell)} - \Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top})\mathcal{Z})^{-1}.$$

Then the integer program (15) is equivalent to the following integer program

$$\min_{Z} \operatorname{tr}(Z^{\top}\Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top}Z)
\text{s.t.} \quad z_{it} \leq z_{i,t+1}
z_{it} \in \{-1,1\}$$
(24)

The covex relaxation of the integer program (24) is provided in the following lemma.

LEMMA 4. Suppose Assumptions 1-2 hold, the potential outcome follows model (14) and does not have observed or latent covariates, i.e., $Y_{it}(z_{it}) = \alpha_i + \beta_t + \tau_1 z_{it} + \tau_2 z_{i,t-1} + \cdots + \tau_{\ell+1} z_{i,t-\ell} + \varepsilon_{it}$, then the convex relaxation of the integer program (24) is

$$\min_{\omega} \sum_{j=1}^{\ell+1} \left[\sum_{t=j}^{T-\ell-1+j} \omega_t^2 - \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell-1+j} \omega_t \right)^2 + \sum_{t=1}^{T-\ell-1+j} \frac{2(T-\ell-1+2j-2t)}{T-\ell} \omega_t \right] \\
s.t. \quad -1 \le \omega_t \le 1 \\
\omega_t \le \omega_{t+1}$$
(25)

The objective function in the quadratic program (25) is a sum of $\ell + 1$ quadratic problems and it is equal to the objective function (19) in the overview of the proof of Theorem 5. If we minimize each one of these $\ell + 1$ problems, the optimal solution of the sub-problem is the same as that in Theorem 1 with the number of time periods to be $T - \ell$. The proof of Lemma 4 and the remaining steps to prove Theorem 5 are provided in Appendix E.

REMARK 16. Beyond the examples for ω^* when $\ell = 1$ and $\ell = 2$ in Remark 8, we also provide the expression for ω^* when $\ell = 3$,

$$\begin{split} &\omega_1^* = -1, \quad \omega_2^* = -1 + \frac{6}{6T^2 - 44T + 79}, \quad \omega_3^* = -1 + \frac{12(T-4)}{6T^2 - 44T + 79}, \\ &\omega_t^* = -1 + \frac{2t-4}{T-3} \quad \text{for } t = 4, \cdots, T-3, \\ &\omega_{T-2}^* = 1 - \frac{12(T-4)}{6T^2 - 44T + 79}, \quad \omega_{T-1}^* = 1 - \frac{6}{6T^2 - 44T + 79}, \quad \omega_T^* = 1 \end{split}$$

If we not only care about the variance of estimated direct and lagged effects, but also the covariance between estimated direct and lagged effects, i.e., $\operatorname{Cov}(\hat{\tau}_j, \hat{\tau}_{j'})$ for $j \neq j'$, then it is natural to consider minimizing $\det(\operatorname{Var}(\hat{\vec{\tau}}))$. Note that when we do not have latent covariates, $\operatorname{Var}(\hat{\vec{\tau}}) = 1/\det(\mathcal{Z}^{\top}(I_{N(T-\ell)} - \Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top})\mathcal{Z})$, minimizing $\operatorname{Var}(\hat{\vec{\tau}})$ is equivalent to minimizing $-\det(\mathcal{Z}^{\top}(I_{N(T-\ell)} - \Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top})\mathcal{Z})$ and the corresponding optimization problem is

$$\min_{\vec{z}=[z_{it}]} -\det(\mathcal{Z}^{\top}(I_{N(T-\ell)} - \Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top})\mathcal{Z})
s.t. z_{it} \leq z_{i,t+1}
z_{it} \in \{-1,1\}$$
(26)

The optimal solution Z for the integer program (26) is called the \mathbf{D} (determinant)-optimal design. We provide the relaxation for the integer program (26) in Lemma 5.

LEMMA 5. Suppose Assumptions 1-2 hold, the potential outcome follows model (14) and does not have observed or latent covariates, i.e., $Y_{it}(z_{it}) = \alpha_i + \beta_t + \tau_1 z_{it} + \tau_2 z_{i,t-1} + \cdots + \tau_{\ell+1} z_{i,t-\ell} + \varepsilon_{it}$, then the relaxation for the integer program (26) is

$$\min_{\vec{\omega}} -det(\Theta)$$

$$s.t. -1 \le \omega_t \le 1$$

$$\omega_t \le \omega_{t+1},$$
(27)

where Θ is defined as

$$\Theta_{jm} = \begin{cases} -N \left[\sum_{t=j}^{T-\ell-1+j} \omega_t^2 - \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell-1+j} \omega_t \right)^2 + \frac{T-\ell}{2} \sum_{t=1}^{T} (v_{T+1-t}^{(j,j)} - v_{T-t}^{(j,j)}) \omega_t \right] & j = m \\ -N \left[\sum_{t=j}^{T-\ell-1+j} \omega_t \omega_{t+m-j} - \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell-1+j} \omega_t \right) \left(\sum_{t=m}^{T-\ell-1+m} \omega_t \right) \right. \\ & \left. + \frac{T-\ell}{2} \sum_{t=1}^{T} (v_{T+1-t}^{(j,m)} - v_{T-t}^{(j,m)}) \omega_t - \sum_{t=j}^{m-1} (\omega_t - \omega_{T-\ell+t}) \right] & j < m \\ \Theta_{mj} & j > m \end{cases}$$

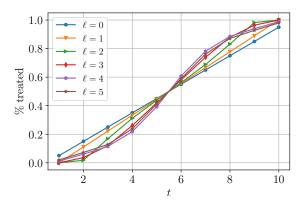


Figure 4 D-optimal treatment design: Optimal treated proportion at each period for a T-period treatment design in the presence of carryover treatment effects, where T = 10. Different colors represent the number of lags ℓ in the carryover effects.

and $v_t^{(j,m)}$ is defined as

$$v_t^{(j,m)} = \begin{cases} 1 & t \leq \ell + 1 - m \\ -\left(-1 + \frac{2(t-1-\ell+m)}{T-\ell}\right) & \ell + 1 - m < t \leq \ell + 1 - j \\ \left(-1 + \frac{2(t-1-\ell+m)}{T-\ell}\right) \left(-1 + \frac{2(t-1-\ell+j)}{T-\ell}\right) & \ell + 1 - j < t \leq T + 1 - m \\ \left(-1 + \frac{2(t-1-\ell+j)}{T-\ell}\right) & T + 1 - m < t \leq T + 1 - j \\ 1 & T + 1 - j < t \end{cases}$$

Note that each entry in Θ is a quadratic function of ω_t , the cross-sectional average of the treatment variables. However, the objective function in the optimization problem (27) is not a convex function of ω_t . In practice, we can use the off-the-shelf software to find the optimal ω_t . We provide the optimal solution for the optimization problem (27) when T = 10 in Figure 4. Similar as all the previous results, the treated proportion is symmetric with respect to the coordinate $(\frac{T+1}{2}, 0.5)$. Similar as the **T**-optimal design in Theorem 5, if the treatment affects more periods (ℓ is larger), the optimal treated proportion is in general smaller at the beginning, increases at a faster rate in the middle, and is larger in the end. The proof of Lemma 5 is provided in Appendix E.

Appendix B: Additional Algorithm

B.1. Grid Search the Tuning Parameter in Low Rank Matrix Estimation with Fixed Effects

When we use LRME to estimate the treatment effect τ given a treatment design Z, we need to specify the regularization parameter μ . If we have the historical control data, we could use *grid search* to find the optimal μ . The details are presented in Algorithm 3. Similar to Section 4.1, we split the historical data into sub-matrices. Given a treatment design matrix Z, we apply the synthetic experiment to each of these sub-matrices and use LRME to estimate τ using a particular μ . The optimal μ minimizes $\max_j |\hat{\tau}^{(j)} - \tau|$ where $\hat{\tau}^{(j)}$ is estimated from the j-th sub-matrix. The next step is to use this treatment design matrix Z to generate synthetic experimental data and estimate τ by LRME using the optimal μ .

Algorithm 3: Grid search optimal μ

```
Inputs: Y^{\text{hist},(1)}, \dots, Y^{\text{hist},(m)}, Z, \tau, \mu_{\text{max}}, \mu_{\text{min}}, N_{\text{grids}}, k_0, \Delta_{\tau}, and t_{\text{max}}
\vec{\mu} \leftarrow \text{logspace}(\mu_{\min}, \mu_{\max}, N_{\text{grids}});
for j = 1, \dots, m do
 \hat{\tau}^{(j)}, \hat{\alpha}^{(j)}, \hat{\beta}^{(j)}, \hat{L}^{(j)} \leftarrow \text{Algorithm } 2(Y^{\text{hist},(j)}, Z, k_0, \mu_{\text{max}}, \Delta_{\tau}, \text{ and } t_{\text{max}});
end
E^{\text{opt}} \leftarrow \max_{i} |\hat{\tau}^{(j)} - \tau|;
\mu^{\text{opt}} = \mu_{\text{max}};
for \mu \in \vec{\mu} do
        for j = 1, \dots, m do
               \hat{\tau}^{(j)}, \hat{\alpha}^{(j)}, \hat{\beta}^{(j)}, \hat{L}^{(j)} \leftarrow \text{Algorithm } 2(Y^{\text{hist},(j)}, Z, k_0, \mu, \Delta_{\tau}, \text{ and } t_{\text{max}});
         end
        \mathbf{E}^{\text{current}} \leftarrow \max_{i} |\hat{\tau}^{(i)} - \tau|;
        \begin{array}{l} \textbf{if} \ E^{current} < E^{opt} \ \textbf{then} \\ \mid \ E^{opt} \leftarrow E^{current}, \ \mu^{opt} \leftarrow \mu; \end{array}
         end
end
Outputs: \mu^{\text{opt}}
```

B.2. Feasible GLS

In feasible GLS, we first estimate τ and residuals from OLS, then estimate errors' covariance matrix from the residuals, and in the end, estimate τ from regression weighted by errors' inverse covariance matrix. We estimate errors' covariance matrix following three assumptions: 1. errors are time-series uncorrelated and identically distributed over time; 2. errors' covariance matrix is the sum of a low-rank matrix and a diagonal matrix (the same assumptions as Theorem 3). We estimate errors' covariance matrix Ω at any time by

$$U, S, U^{\top} = \text{SVD}\left(\frac{1}{T}\hat{e}\hat{e}^{\top}\right), \quad \hat{\Omega} = U_{:,1:k_0}S_{1:k_0,1:k_0}U_{:,1:k_0}^{\top} + \text{diag}\left(\frac{1}{T}\hat{e}\hat{e}^{\top} - U_{:,1:k_0}S_{1:k_0,1:k_0}U_{:,1:k_0}^{\top}\right),$$

where S is a diagonal matrix with diagonal entries to be the singular values of $\frac{1}{T}\hat{e}\hat{e}^{\top}$ in descending order, diag $\left(\frac{1}{T}\hat{e}\hat{e}^{\top} - U_{:,1:k_0}S_{1:k_0,1:k_0}U_{:,1:k_0}^{\top}\right)$ is a diagonal matrix with diagonal entries to be the diagonal entries of $\frac{1}{T}\hat{e}\hat{e}^{T} - U_{:,1:k_0}S_{1:k_0,1:k_0}U_{:,1:k_0}^{\top}$ and $k_0 < T$. The full estimated errors' covariance matrix $(NT \times NT)$ is a block diagonal matrix with each block to be $\hat{\Omega}$. By assuming errors' are time-series uncorrelated, we estimate much fewer parameters and each parameter is estimated from T observations.

Appendix C: No Covariates

Proof of Equation (7). We use linear regression to estimate τ and η , so

$$\begin{bmatrix} \hat{\tau} \\ \hat{\eta} \end{bmatrix} = \left(\begin{bmatrix} \vec{z}^\top \\ \Gamma^\top \end{bmatrix} \begin{bmatrix} \vec{z} \ \Gamma \end{bmatrix} \right)^{-1} \begin{bmatrix} \vec{z}^\top \\ \Gamma^\top \end{bmatrix} \vec{y}$$

 $^{^9} k_0$ can be a preset value or can be chosen in a data-driven manner.

$$= \begin{bmatrix} \tau \\ \eta \end{bmatrix} + \left(\begin{bmatrix} \vec{z}^\intercal \vec{z} \ \vec{z}^\intercal \Gamma \\ \Gamma^\intercal \vec{z} \ \Gamma^\intercal \Gamma \end{bmatrix} \right)^{-1} \begin{bmatrix} \vec{z}^\intercal \\ \Gamma^\intercal \end{bmatrix} \vec{\varepsilon}$$

Since $Var(\vec{\varepsilon}) = \sigma^2 I$,

$$\operatorname{Var}\left(\begin{bmatrix}\hat{\tau}\\\hat{\eta}\end{bmatrix}\right) = \sigma^2 \left(\begin{bmatrix}\vec{z}^\top \vec{z} \ \vec{z}^\top \Gamma\\\Gamma^\top \vec{z} \ \Gamma^\top \Gamma\end{bmatrix}\right)^{-1}$$

From block matrix inversion, we have

$$\operatorname{Var}(\hat{\tau}) = \sigma^2 \left(\vec{z}^\top \vec{z} - \vec{z}^\top \Gamma (\Gamma^\top \Gamma)^{-1} \Gamma^\top \vec{z} \right)^{-1}$$

which is equivalent to Eq. (7). \square

Proof of Lemma 2. When there are only time fixed effects, Γ can be simplified to

$$\Gamma = I_T \otimes \vec{1} = \begin{bmatrix} \vec{1} & & \\ & \vec{1} & \\ & \ddots & \\ & & \vec{1} \end{bmatrix} \in \mathbb{R}^{NT \times T}, \quad \text{ where } \vec{1} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^N$$

Then we have $\Gamma^{\top}\Gamma = N \cdot I_T$. $\frac{1}{N}\vec{z}^{\top}\Gamma = \left[\omega_1 \cdots \omega_T\right]$, where $\omega_t = \frac{1}{N}\sum_{i=1}^N z_{it}$. The objective function in the integer program (8) has

$$\vec{z}^{\mathsf{T}} \Gamma (\Gamma^{\mathsf{T}} \Gamma)^{-1} \Gamma^{\mathsf{T}} \vec{z} = N \sum_{t=1}^{T} \omega_t^2$$

It is minimized at $\omega_t = 0$. Thus, any treatment design is optimal if it satisfies $\omega_t = 0$ for all t and $z_{it} = z_{i,t+1}$ for all i and t.

In other words, if N is even, $\operatorname{Var}(\hat{\tau})$ is minimized at $\omega_t = 0$, implying 50% of z_{it} to be 1 and the other 50% of z_{it} to be -1. If N is odd, $\operatorname{Var}(\hat{\tau})$ is minimized at $|\omega_t| = \frac{1}{N}$, implying there are $\frac{N-1}{2}$ units with $z_{it} = 1$ for all t, there are $\frac{N-1}{2}$ units with $z_{it} = -1$ for all t, and there is one unit with z_{it} to be flexible as long as $z_{it} \leq z_{i,t+1}$ is satisfied. \square

Proof of Lemma 3. When there are only unit fixed effects, Γ can be simplified to

$$\Gamma = \vec{1} \otimes I_N = \begin{bmatrix} I_N \\ I_N \\ \vdots \\ I_N \end{bmatrix} \in \mathbb{R}^{NT \times N}, \quad \text{ where } \vec{1} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^T.$$

¹⁰ Then we have $\Gamma^{\top}\Gamma = T \cdot I_N$. $\frac{1}{T}\vec{z}^{\top}\Gamma = \left[\zeta_1 \cdots \zeta_N\right]$, where $\zeta_i = \frac{1}{T}\sum_{t=1}^T z_{it}$. The objective function in the integer program (8) has

$$\vec{z}^{\top} \Gamma (\Gamma^{\top} \Gamma)^{-1} \Gamma^{\top} \vec{z} = T \sum_{i=1}^{N} \zeta_i^2.$$

It is minimized at $\omega_t = 0$. Thus, any treatment design is optimal if it satisfies $\zeta_i = 0$ for all i. Given Assumption 2, the equivalent statement is that any treatment design is optimal if it satisfies $\omega_t = -1$ for all $t < \frac{T+1}{2}$ and $\omega_t = 1$ for all $t > \frac{T+1}{2}$.

In other words, if T is even, $\operatorname{Var}(\hat{\tau})$ is minimized at $\zeta_i = 0$ for all i, implying all units are untreated in the first half time periods and they all switch to the treated group in the second half time periods. If T is odd, $\operatorname{Var}(\hat{\tau})$ is minimized at $|\zeta_i| = \frac{1}{T}$, implying all units are in the control group for all $t \leq \frac{T-1}{2}$, they can be either treated or untreated at time period $t = \frac{T+1}{2}$, and they are in the treated group for all $t > \frac{T+1}{2}$. \square

¹⁰ Since we not have time fixed effects, we do not have the identification problem and we do not need to assume $\alpha_N = 0$ as a contrast to Theorem 1.

Proof of Theorem 1. Denote $\zeta_i = \frac{1}{T} \sum_{t=1}^T z_{it}$ for all i and $\omega_t = \frac{1}{N} \sum_{i=1}^N z_{it}$ for all t.

There are two steps to show Theorem 1:

- 1. The objective function in the integer program (8), $\vec{z}^{\top}\Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top}\vec{z}$, can be simplified to $\vec{\omega}^{\top}P_{\vec{1}}\vec{\omega} + 2\vec{d}^{\top}\vec{\omega}$ with $P_{\vec{1}} = I_T \frac{1}{T}\vec{1}\vec{1}^{\top} \in \mathbb{R}^{T \times T}$, $\vec{d} = [d_t] \in \mathbb{R}^T$, and $d_t = \frac{T+1-2t}{T}$ for all t.
 - 2. $\omega_t^* = \frac{2t-1-T}{T}$ is the optimal solution that satisfies the first order condition (FOC) of $\vec{\omega}^{\top} P_{\vec{1}} \vec{\omega} + 2 \vec{d}^{\top} \vec{\omega}$. Now we start with the first step. There are two main intermediate steps.
- 1. Express $\vec{z}^{\top}\Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top}\vec{z}$ in terms of ω_t and ζ_i . When potential outcomes follow $Y_{it}(z_{it}) = \alpha_i + \beta_t + \tau z_{it} + \varepsilon_{it}$, we have

$$\Gamma = \begin{bmatrix} \Gamma_{:1} \\ \Gamma_{:2} \\ \vdots \\ \Gamma_{:T} \end{bmatrix} = \begin{bmatrix} \tilde{I} & \tilde{I} \\ \tilde{I} & \tilde{I} \\ \vdots & \ddots \\ \tilde{I} & & \tilde{I} \end{bmatrix} \in \mathbb{R}^{NT \times (N+T-1)}.$$

where $\tilde{I} = \begin{bmatrix} I_{N-1} & \vec{0} \end{bmatrix}^\top \in \mathbb{R}^{N \times (N-1)}, \ \vec{0} \in \{0\}^{N-1} \text{ and } \vec{1} \in \{1\}^N. \text{ In } z^\top \Gamma(\Gamma^\top \Gamma)^{-1} \Gamma^\top z, \text{ we have } \vec{0} = \vec{0}$

$$(\Gamma^{\top}\Gamma)^{-1} = \begin{bmatrix} \Xi_{11} \ \Xi_{12} \\ \Xi_{21} \ \Xi_{22} \end{bmatrix} = \begin{bmatrix} M & -\frac{1}{N}MM_{\vec{1}} \\ -\frac{1}{N}M_{\vec{1}}^{\top}M \ \frac{1}{N}I + \frac{1}{N^2}M_{\vec{1}}^{\top}MM_{\vec{1}} \end{bmatrix} \in \mathbb{R}^{(N+T-1)\times(N+T-1)},$$

where $\Xi_{11} = M$, $\Xi_{12} = -\frac{1}{N} M M_{\vec{1}}$, $\Xi_{21} = -\frac{1}{N} M_{\vec{1}}^{\top} M$, $\Xi_{22} = \frac{1}{N} I + \frac{1}{N^2} M_{\vec{1}}^{\top} M M_{\vec{1}}$ and

$$M = \frac{1}{T} \left(I_{N-1} - \frac{1}{N} \vec{1} \vec{1}^{\top} \right)^{-1} = \frac{1}{T} (I + \vec{1} \vec{1}^{\top}) \in \mathbb{R}^{(N-1) \times (N-1)}, \quad M_{\vec{1}} = \left[\vec{1} \cdots \vec{1} \right] \in \mathbb{R}^{(N-1) \times T}.$$

where $\vec{1} = \{1\}^{N-1}$ in M and $M_{\vec{1}}$. Furthermore,

$$\vec{z}^{\top} \Gamma = \left[\sum_{t=1}^{T} z_{1t} \, \cdots \, \sum_{t=1}^{T} z_{N-1,t} \, \sum_{i=1}^{N} z_{i1} \, \cdots \, \sum_{i=1}^{N} z_{iT} \right] = \left[T \zeta_{1} \, \cdots \, T \zeta_{N-1} \, N \omega_{1} \, \cdots \, N \omega_{T} \right],$$

Let $\phi \coloneqq \begin{bmatrix} T\zeta_1 & \cdots & T\zeta_{N-1} \end{bmatrix}^{\top}$ and $\iota \coloneqq \begin{bmatrix} N\omega_1 & \cdots & N\omega_T \end{bmatrix}^{\top}$. Then we have $\vec{z}^{\top}\Gamma = \begin{bmatrix} \phi^{\top} & \iota^{\top} \end{bmatrix}$ and

$$\vec{z}^\top \Gamma (\Gamma^\top \Gamma)^{-1} \Gamma^\top \vec{z} = \phi^\top \Xi_{11} \phi + 2 \phi^\top \Xi_{12} \iota + \iota^\top \Xi_{22} \iota.$$

We simplify each term in $\vec{z}^{\top}\Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top}\vec{z}$,

$$\begin{split} \phi^{\top}\Xi_{11}\phi &= \frac{1}{T}\phi^{\top}(I_{N-1} + \vec{1}\vec{1}^{\top})\phi = T\sum_{i=1}^{N-1}\zeta_{i}^{2} + T\left(\sum_{i=1}^{N-1}\zeta_{i}\right)^{2} \\ \phi^{\top}\Xi_{12}\iota &= -\frac{1}{T}\phi^{\top}(I_{N-1} + \vec{1}\vec{1}^{\top})\vec{1}\left(\sum_{t=1}^{T}\omega_{t}\right) = -\frac{N}{T}\phi^{\top}\vec{1}\left(\sum_{t=1}^{T}\omega_{t}\right) = -N\left(\sum_{i=1}^{N-1}\zeta_{i}\right)\left(\sum_{t=1}^{T}\omega_{t}\right) \\ \iota^{\top}\Xi_{22}\iota &= N\sum_{t=1}^{T}\omega_{t}^{2} + \frac{N(N-1)}{T}\left(\sum_{t=1}^{T}\omega_{t}\right)^{2} \end{split}$$

From the definitions of ζ_i and ω_t , $T\sum_{i=1}^N \zeta_i = N\sum_{t=1}^T \omega_t$. Then $\sum_{i=1}^{N-1} \zeta_i = \frac{N}{T}\sum_{t=1}^T \omega_t - \zeta_N$. Using this property, we have

$$\vec{z}^{\top} \Gamma (\Gamma^{\top} \Gamma)^{-1} \Gamma^{\top} \vec{z} = T \sum_{i=1}^{N} \zeta_i^2 - \frac{N}{T} \left(\sum_{t=1}^{T} \omega_t \right)^2 + N \sum_{t=1}^{T} \omega_t^2.$$
 (28)

2. Express $\vec{z}^{\top}\Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top}\vec{z}$ in terms of ω_t . Following Assumption 2, $\sum_{i=1}^{N} \zeta_i^2$ is fixed given ω_t . Given ω_t , there are $\frac{N(1+\omega_1)}{2}, \frac{N(1+\omega_2)}{2}, \cdots \frac{N(1+\omega_T)}{2}$ treated units in time period $1, 2, \cdots, T$. It is equivalent to having $\frac{N(1+\omega_1)}{2}, \frac{N(\omega_2-\omega_1)}{2}, \cdots, \frac{N(\omega_T-\omega_{T-1})}{2}$ untreated units to start the treatment at time period $1, 2, \cdots, T$ and leaving $\frac{N(1-\omega_T)}{2}$ units in the control group in the end. Then we have

$$\sum_{i=1}^{N} \zeta_i^2 = N \left[\frac{1+\omega_1}{2} \cdot 1 + \frac{\omega_2 - \omega_1}{2} \left(\frac{T-2}{T} \right)^2 + \frac{\omega_3 - \omega_2}{2} \left(\frac{T-4}{T} \right)^2 + \cdots \right]$$

$$+ \frac{\omega_T - \omega_{T-1}}{2} \left(-1 + \frac{2}{T} \right)^2 + \frac{1-\omega_T}{2} \cdot 1$$

$$= N \left[1 + \left(2 - \frac{2}{T} \right) \frac{\omega_1}{T} + \left(2 - \frac{6}{T} \right) \frac{\omega_2}{T} + \cdots + \left(2 - \frac{4T-2}{T} \right) \frac{\omega_T}{T} \right]$$

$$= \frac{2N}{T} \sum_{t=1}^{T} \frac{T+1-2t}{T} \cdot \omega_t$$

Therefore, we can write $\vec{z}^{\top}\Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top}\vec{z}$ as a function of ω_t and denote $\vec{z}^{\top}\Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top}\vec{z}$ as $f(\vec{\omega})$

$$f(\vec{\omega}) \coloneqq \vec{z}^{\top} \Gamma(\Gamma^{\top} \Gamma)^{-1} \Gamma^{\top} \vec{z} = N \sum_{t=1}^{T} \omega_t^2 - \frac{N}{T} \left(\sum_{t=1}^{T} \omega_t \right)^2 + \frac{2N}{T} \sum_{t=1}^{T} (T+1-2t) \cdot \omega_t = N \left(\vec{\omega}^{\top} P_{\vec{1}} \vec{\omega} + 2 \vec{d}^{\top} \vec{\omega} \right).$$

Next we show the second step. Let $f_1(\vec{\omega}) := \sum_{t=1}^T \omega_t^2$, $f_2(\vec{\omega}) := -\frac{1}{T} \left(\sum_{t=1}^T \omega_t\right)^2$ and $f_3(\vec{\omega}) := 2\sum_{t=1}^T \frac{T+1-2t}{T} \omega_t$. Then we have $f(\vec{\omega}) = f_1(\vec{\omega}) + f_2(\vec{\omega}) + f_3(\vec{\omega})$. We first take the second order derivative

$$\nabla^2 f_1(\vec{\omega}) = 2I_T$$

$$\nabla^2 f_2(\vec{\omega}) = -\frac{2}{T} \vec{1} \vec{1}^{\top}$$

$$\nabla^2 f_3(\vec{\omega}) = 0$$

Thus,

$$\nabla^2 f(\vec{\omega}) = \nabla^2 f_1(\vec{\omega}) + \nabla^2 f_2(\vec{\omega}) + \nabla^2 f_3(\vec{\omega}) \ge 0$$

and $f(\vec{\omega})$ is convex. Next we take the first order derivative

$$\nabla f(\vec{\omega}) = \left[2\omega_1 + \frac{2(T-1)}{T} - \frac{2}{T}\sum_{t=1}^{T}\omega_t \ 2\omega_2 + \frac{2(T-3)}{T} - \frac{2}{T}\sum_{t=1}^{T}\omega_t \ \cdots \ 2\omega_T - \frac{2(T-1)}{T} - \frac{2}{T}\sum_{t=1}^{T}\omega_t\right]$$

When $\omega_t^* = \frac{2t-1-T}{T}$, $\nabla f(\vec{\omega}^*) = 0$. Thus, $\vec{\omega}^*$ is a local optimal solution. Since $f(\vec{\omega})$ is convex, the local minimum is also the global minimum. Therefore, $\vec{\omega}^*$ is the global optimal solution.

Proof of Theorem 2 Denote the objective function of the quadratic program (9) as $f(\vec{\omega})$ (similar as the proof of Theorem 1). The variance of τ (denoted as $g(\vec{\omega})$) is

$$g(\vec{\omega}) \coloneqq \operatorname{Var}(\hat{\tau}) = \frac{\sigma^2}{\vec{z}^\top (I_{NT} - \Gamma(\Gamma^\top \Gamma)^{-1} \Gamma^\top) \vec{z}} = \frac{\sigma^2}{N(T - f(\vec{\omega}))}$$

The optimal solution in Theorem 1 is $\omega_t^* = \frac{2t-1-T}{T}$. We plug $\vec{\omega}^* = [\omega_t^*]$ into $f(\vec{\omega})$,

$$f(\vec{\omega}^*) = -\sum_{t=1}^{T} \left(\frac{T+1-2t}{T}\right)^2 = -\frac{(T+1)(T-1)}{3T}.$$

The rounded solution based on $\vec{\omega}^*$ is denoted as $\vec{\omega}^{\text{rnd}} = [\omega_t^{\text{rnd}}]$ (from the nearest integer rounding rule). We have $|\omega_t^{\text{rnd}} - \omega_t^*| \leq \frac{1}{N}$ and $|\sum_{t=1}^T \omega_t^{\text{rnd}}| \leq \frac{1}{N}$. Therefore,

$$f(\vec{\omega}^{\text{rnd}}) = \sum_{t=1}^{T} \left(\omega_t^{\text{rnd}} + \frac{T+1-2t}{T} \right)^2 - \frac{1}{T} \left(\sum_{t=1}^{T} \omega_t^{\text{rnd}} \right)^2 - \sum_{t=1}^{T} \left(\frac{T+1-2t}{T} \right)^2 \\ \leq \frac{T}{N^2} - \frac{(T+1)(T-1)}{3T}$$

We have

$$\frac{g(\vec{\omega}^{\text{rnd}})}{g(\vec{\omega}^*)} = \frac{T - f(\vec{\omega}^*)}{T - f(\vec{\omega}^{\text{rnd}})} \leq \frac{(4T^2 - 1)/(3T)}{(4T^2 - 1)/(3T) - T/N^2} = \frac{4T^2 - 1}{4T^2 - 1 - 3T^2/N^2} \leq \frac{1}{1 - 1/N^2}$$

following $3T^2 \leq 4T^2 - 1$ for all T. Note that $f(\vec{\omega}^*)$ reaches the minimum objective function value in the quadratic program (9), so $f(\vec{\omega}^*)$ is the lower bound for the minimum objective function value in the original integer programming problem. Thus, $g(\vec{\omega}^*) \leq \operatorname{Var}_{Z^*}(\hat{\tau})$ and we have

$$\operatorname{Var}_{Z^{\operatorname{rnd}}}(\hat{\tau}) \leq \frac{1}{1 - 1/N^2} \operatorname{Var}_{Z^*}(\hat{\tau}).$$

Appendix D: With Covariates

Proof of Theorem 3. We can combine β_t with $X_i^{\top}\theta_t$ in the potential outcome model (11), that is,

$$Y_{it}(z_{it}) = \alpha_i + \underbrace{\begin{bmatrix} 1 \ X_i^{\top} \end{bmatrix}}_{\tilde{X}_i^{\top}} \begin{bmatrix} \beta_t \\ \theta_t \end{bmatrix} + \tau z_{it} + \varepsilon_{it},$$

and denote p := r + 1 so $\tilde{X}_i \in \mathbb{R}^p$. Denote $\zeta_i = \frac{1}{T} \sum_{t=1}^T z_{it}$ for all i and $\tilde{\omega}_t = \frac{1}{N} \sum_{i=1}^N \tilde{X}_i z_{it} \in \mathbb{R}^p$ for all t. First, when we use $W = \sum_{e}^{-1}$ in GLS, with some algebra, we have the following expression for the variance of $\hat{\tau}$,

$$\operatorname{Var}(\hat{\tau}) = \vec{z}^{\top} \left(\Sigma_e^{-1} - \Sigma_e^{-1} \Gamma (\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1} \right) \vec{z}, \tag{29}$$

where $\tilde{I} = \begin{bmatrix} I_{N-p} \ \mathbf{0}_{N-p,p} \end{bmatrix}^{\top} \in \mathbb{R}^{N \times (N-p)}$, $\mathbf{0}_{N-p,p}$ is a matrix of 0, $\vec{1} \in \{1\}^N$, and

$$\Gamma = \begin{bmatrix} \Gamma_{:1} \\ \Gamma_{:2} \\ \vdots \\ \Gamma_{:T} \end{bmatrix} = \begin{bmatrix} \tilde{I} & \tilde{I} & X \\ \tilde{I} & \tilde{I} & X \\ \vdots & & \ddots \\ \tilde{I} & & & \tilde{I} & X \end{bmatrix} = \begin{bmatrix} \tilde{I} & \tilde{X} \\ \tilde{I} & \tilde{X} \\ \vdots & & \ddots \\ \tilde{I} & & & \tilde{X} \end{bmatrix} \in \mathbb{R}^{NT \times (N-p+Tp)}.$$
(30)

We restrict $\alpha_i = 0$ for $i = N - p, \dots, N$ such that all other α_i and β_t can be uniquely identified. Details to show Eq. (29) are the same as those to show Eq. (7).

There are three steps to show Theorem 3:

1. Minimizing $Var(\hat{\tau})$ can be simplified to minimize

$$\begin{split} & (\vec{\omega}^{(1)})^{\top} P_{\vec{1}} \vec{\omega}^{(1)} + 2 \vec{d}^{\top} \vec{\omega}^{(1)} + \sum_{j=2}^{p} (\vec{\omega}^{(j)})^{\top} P_{\vec{1}} \vec{\omega}^{(j)} \\ & + \frac{1}{N} \sum_{t=1}^{T} \vec{z}_{t}^{\top} U(I_{k} + U^{\top} U)^{-1} U^{\top} \vec{z}_{t} - \frac{1}{NT} \left(\sum_{t=1}^{T} \vec{z}_{t} \right) U(I_{k} + U^{\top} U)^{-1} U^{\top} \left(\sum_{t=1}^{T} \vec{z}_{t} \right), \end{split}$$

where $\vec{\omega}^{(j)} = \begin{bmatrix} \tilde{\omega}_{1,j} & \cdots & \tilde{\omega}_{T,j} \end{bmatrix}^{\top}$ and $\tilde{\omega}_{t,j} = \frac{1}{N} \sum_{i=1}^{N} \tilde{X}_{i,j} z_{it}$ for $j = 1, \cdots, p$ and $t = 1, \cdots, T$.

2. Take the first-order condition and show sufficient conditions for optimality.

- 3. Based on the sufficient conditions and finite u_i , show $\omega_{g,t} = \frac{2t-1-T}{T}$ is optimal. We start with the first step. The steps to simplify $\vec{z}^{\top} \left(\Sigma_e^{-1} - \Sigma_e^{-1} \Gamma(\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1} \right) \vec{z}$ are as follows.
- 1. Calculate Σ_e^{-1} . Under the assumptions that v_t has mean 0 and variance I_k , $Cov(v_t, v_{t-q}) = 0$ for any q, $Cov(v_t, \varepsilon_{is}) = 0$ for any s, t and i, and $\sigma = 1$, we have $\Sigma_e^{-1} = \text{diag}(\Psi, \Psi, \dots, \Psi)$, where

$$\Psi = \Sigma_{e_{\star}}^{-1} = (I_N + UU^{\top})^{-1} = I_N - U(I_k + U^{\top}U)^{-1}U^{\top}$$

- 2. Calculate $\vec{z}^{\top} \Sigma_e^{-1} \vec{z}$. We have $z^{\top} \Sigma_e^{-1} z = \sum_{t=1}^T z_t^{\top} \Psi z_t$.
- 3. Calculate $\vec{z}^{\top} \Sigma_e^{-1} \Gamma$. We have

$$\vec{z}^{\top} \Sigma_e^{-1} \Gamma = \vec{z}^{\top} \begin{bmatrix} \Psi \tilde{I} & \Psi \tilde{X} \\ \Psi \tilde{I} & \Psi \tilde{X} \\ \vdots & & \ddots \\ \Psi \tilde{I} & & \Psi \tilde{X} \end{bmatrix} = \left[\sum_{t=1}^T \vec{z}_t^{\top} \Psi \tilde{I} & \vec{z}_1^{\top} \Psi \tilde{X} & \cdots & \vec{z}_T^{\top} \Psi \tilde{X} \right] = \left[\phi^{\top} \ \iota^{\top} \right],$$

where $\phi^{\top} = \sum_{t=1}^{T} \vec{z}_{t}^{\top} \Psi \tilde{I} = (\sum_{t=1}^{T} \vec{z}_{t})^{\top} \Psi \tilde{I}$ and $\iota^{\top} = \left[\vec{z}_{1}^{\top} \Psi \tilde{X} \cdots \vec{z}_{T}^{\top} \Psi \tilde{X} \right]$. Let $\zeta \coloneqq \frac{1}{T} \sum_{t=1}^{T} \vec{z}_{t} \in \mathbb{R}^{N}$. Then $\phi = T \zeta^{\top} \Psi \tilde{I}$. Note that U and \tilde{X} are orthogonal, we have $\Psi \tilde{X} = \tilde{X}$ and $\tilde{X}^{\top} \Psi \tilde{X} = N \cdot I_{p}$. Then $\iota^{\top} = \left[\vec{z}_{1}^{\top} \tilde{X} \cdots \vec{z}_{T}^{\top} \tilde{X} \right] = \left[N \tilde{\omega}_{1}^{\top} \cdots N \tilde{\omega}_{T}^{\top} \right] = N \tilde{\omega}^{\top} \in \mathbb{R}^{T_{p}}$, where $\tilde{\omega}_{t} = \frac{1}{N} \sum_{i=1}^{N} \tilde{X}_{i} z_{it} \in \mathbb{R}^{p}$.

4. Calculate $(\Gamma^{\top}\Sigma_e\Gamma)^{-1}$. Note that

$$(\Gamma^{\top} \Sigma_e \Gamma)^{-1} = \begin{bmatrix} \Xi_{11} & \Xi_{12} \\ \Xi_{21} & \Xi_{22} \end{bmatrix} = \begin{bmatrix} M & -M\tilde{M} \\ -\tilde{M}^{\top} M & \bar{M} + \tilde{M}^{\top} M \tilde{M} \end{bmatrix} \in \mathbb{R}^{(N+(T-1)p)) \times (N+(T-1)p)},$$

where $\Xi_{11} = M, \; \Xi_{12} = -M\tilde{M}, \; \Xi_{21} = -\tilde{M}^{\top}M, \; \Xi_{22} = \bar{M} + \tilde{M}^{\top}M\tilde{M} \;$ with

$$M = \frac{1}{T} \left(\tilde{I}^{\top} \Psi \tilde{I} - \tilde{I}^{\top} \Psi \tilde{X} (\tilde{X}^{\top} \Psi \tilde{X})^{-1} \tilde{X}^{\top} \Psi \tilde{I} \right)^{-1} = \frac{1}{T} \left(\tilde{I}^{\top} \Psi \tilde{I} - \frac{1}{N} \tilde{X} \tilde{X}^{\top} \right)^{-1} \in \mathbb{R}^{(N-p) \times (N-p)}$$

$$\tilde{M} = \left[\tilde{I}^{\top} \Psi \tilde{X} (\tilde{X}^{\top} \Psi \tilde{X})^{-1} \cdots \tilde{I}^{\top} \Psi \tilde{X} (\tilde{X}^{\top} \Psi \tilde{X})^{-1} \right] = \left[\frac{1}{N} \tilde{X} \cdots \frac{1}{N} \tilde{X} \right] \in \mathbb{R}^{(N-p) \times Tp}$$

$$\bar{M} = \operatorname{diag}((\tilde{X}^{\top} \Psi \tilde{X})^{-1}, (\tilde{X}^{\top} \Psi \tilde{X})^{-1}, \cdots, (\tilde{X}^{\top} \Psi \tilde{X})^{-1}) = \frac{1}{N} I_{Tp} \in \mathbb{R}^{Tp \times Tp}$$

We can simplify $M = \frac{1}{T} \left(\tilde{I}^\top \Psi \tilde{I} - \frac{1}{N} \tilde{X} \tilde{X}^\top \right)^{-1}$ by

$$M = \frac{1}{T} \left[\left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} + \left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} \tilde{X} \left(N - \tilde{X}^{\top} \left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} \tilde{X} \right)^{-1} \tilde{X}^{\top} \left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} \right]$$

$$\left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} = \left(I_{N-p} - U_{(1)} (I_k + U^{\top} U)^{-1} U_{(1)}^{\top} \right)^{-1} = I_{N-p} + U_{(1)} (I_k + U_{(2)}^{\top} U_{(2)})^{-1} U_{(1)}^{\top}$$

where $U_{(1)} = \begin{bmatrix} u_1 \ u_2 \cdots u_{N-p} \end{bmatrix}^{\top} \in \mathbb{R}^{(N-p)\times k}$ and $U_{(2)} = \begin{bmatrix} u_{N-p+1} \cdots u_N \end{bmatrix}^{\top} \in \mathbb{R}^{p\times k}$ (note $U = \begin{bmatrix} U_{(1)}^{\top} \ U_{(2)}^{\top} \end{bmatrix}^{\top}$). 5. Calculate $\vec{z}^{\top} \Sigma_e^{-1} \Gamma (\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1} \vec{z}$. It is equivalent to calculating $\begin{bmatrix} \phi^{\top} \ \iota^{\top} \end{bmatrix} \begin{bmatrix} \Xi_{11} \ \Xi_{12} \\ \Xi_{21} \ \Xi_{22} \end{bmatrix} \begin{bmatrix} \phi \\ \iota \end{bmatrix} = 0$

5. Calculate $\vec{z}^{\top} \Sigma_e^{-1} \Gamma(\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1} \vec{z}$. It is equivalent to calculating $\begin{bmatrix} \phi^{\top} \ \iota^{\top} \end{bmatrix} \begin{bmatrix} \Xi_{11} & \Xi_{12} \\ \Xi_{21} & \Xi_{22} \end{bmatrix}$ $\phi^{\top} \Xi_{11} \phi + \phi^{\top} \Xi_{12} \iota + \iota^{\top} \Xi_{21} \phi + \iota^{\top} \Xi_{22} \iota$, where (recall $\zeta = \frac{1}{T} \sum_{t=1}^{T} \vec{z}_t$ and $\tilde{\vec{\omega}} = \begin{bmatrix} \tilde{\omega}_1^{\top} & \cdots \tilde{\omega}_T^{\top} \end{bmatrix}^{\top}$)

$$\begin{split} \phi^\top \Xi_{11} \phi &= T^2 \zeta^\top \Psi \tilde{I} M \tilde{I}^\top \Psi \zeta \\ \phi^\top \Xi_{12} \iota &= -N T \zeta^\top \Psi \tilde{I} M \tilde{M} \tilde{\omega} = -T \zeta^\top \Psi \tilde{I} M \tilde{X} \left(\sum_{t=1}^T \tilde{\omega}_t \right) \\ \iota^\top \Xi_{21} \phi &= -N T \tilde{\omega}^\top \tilde{M}^\top M \Psi \tilde{I} \zeta = -T \left(\sum_{t=1}^T \tilde{\omega}_t^\top \right) \tilde{X}^\top M \Psi \tilde{I}^\top \zeta \\ \iota^\top \Xi_{22} \iota &= N \sum_{t=1}^T \tilde{\omega}_t^\top \tilde{\omega}_t + \left(\sum_{t=1}^T \tilde{\omega}_t^\top \right) \tilde{X}^\top M \tilde{X} \left(\sum_{t=1}^T \tilde{\omega}_t \right) \end{split}$$

Let $\tilde{X}_{(1)} \coloneqq \begin{bmatrix} \tilde{X}_1 \ \tilde{X}_2 \ \cdots \ \tilde{X}_{N-p} \end{bmatrix}^\top \in \mathbb{R}^{(N-p)\times p}$ and $\tilde{X}_{(2)} \coloneqq \begin{bmatrix} \tilde{X}_{N-p+1} \ \cdots \ \tilde{X}_N \end{bmatrix}^\top \in \mathbb{R}^{p\times p}$ (note $\tilde{X} = \begin{bmatrix} \tilde{X}_{(1)}^\top \ \tilde{X}_{(2)}^\top \end{bmatrix}^\top$).

We can simplify
$$\phi^{\top}\Xi_{11}\phi$$
, $\phi^{\top}\Xi_{12}\iota$, $\iota^{\top}\Xi_{21}\phi$ and $\iota^{\top}\Xi_{22}\iota$ by calculating the following terms

$$G = \tilde{I} \left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} \tilde{I}^{\top} = \begin{bmatrix} I_{N-p} + U_{(1)} (I_k + U_{(2)}^{\top} U_{(2)})^{-1} U_{(1)}^{\top} & \vec{0} \\ \vec{0}^{\top} \end{bmatrix} \in \mathbb{R}^{N \times N}$$

$$\Omega = \Psi G \Psi = \begin{bmatrix} I_{N-p} & \vec{0} \\ -U_{(2)} (I_k + U_{(2)}^{\top} U_{(2)})^{-1} U_{(1)}^{\top} & 0 \end{bmatrix} \Psi$$

$$= \begin{bmatrix} I_{N-p} - U_{(1)} (I_k + U^{\top} U)^{-1} U_{(1)}^{\top} & -U_{(1)} (I_k + U^{\top} U)^{-1} U_{(2)}^{\top} \\ -U_{(2)} (I_k + U^{\top} U)^{-1} U_{(1)}^{\top} & U_{(2)} (I_k + U_{(2)}^{\top} U_{(2)})^{-1} U_{(1)}^{\top} U_{(1)} (I_k + U^{\top} U)^{-1} U_{(2)}^{\top} \end{bmatrix} \in \mathbb{R}^{N \times N}$$

$$\left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} \tilde{X}_{(1)} = \tilde{X}_{(1)} - U_{(1)} (I_k + U_{(2)}^{\top} U_{(2)})^{-1} U_{(2)}^{\top} \tilde{X}_{(2)} \in \mathbb{R}^{(N-p) \times p}$$

$$\tilde{I} \left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} \tilde{X}_{(1)} = \begin{bmatrix} \tilde{X}_{(1)} - U_{(1)} (I_k + U_{(2)}^{\top} U_{(2)})^{-1} U_{(2)}^{\top} \tilde{X}_{(2)} \\ 0 \end{bmatrix} \in \mathbb{R}^{N \times p}$$

$$\delta = \tilde{X}_{(1)}^{\top} \left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} \tilde{X}_{(1)} = N I_p - (\tilde{X}_{(2)}^{\top} \tilde{X}_{(2)} - \tilde{X}_{(2)}^{\top} U_{(2)} (I_k + U_{(2)}^{\top} U_{(2)})^{-1} U_{(2)}^{\top} \tilde{X}_{(2)} \right) \in \mathbb{R}^{p \times p}$$

$$\gamma = \Psi \tilde{I} \left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} \tilde{X}_{(1)} = \begin{bmatrix} \tilde{X}_{(1)} \\ U_{(2)} (I_k + U_{(2)}^{\top} U_{(2)})^{-1} U_{(2)}^{\top} \tilde{X}_{(2)} \end{bmatrix} \in \mathbb{R}^{N \times p}$$

Together with $N \sum_{t=1}^{T} \tilde{\omega}_t = T \sum_{i=1}^{N} \tilde{X}_i \zeta_i$ and $N \vec{1}^{\top} \omega = T \tilde{X}^{\top} \zeta$, we have

$$\begin{split} \phi^{\top}\Xi_{11}\phi &= T\zeta^{\top}\left(\Omega + \gamma(NI_{p} - \boldsymbol{\delta})^{-1}\gamma^{\top}\right)\zeta\\ \phi^{\top}\Xi_{12}\iota &= -T\zeta^{\top}\left(\gamma(NI_{p} - \boldsymbol{\delta})^{-1}\tilde{X}^{\top}\right)\zeta\\ \iota^{\top}\Xi_{21}\phi &= -T\zeta^{\top}\left(\tilde{X}(NI_{p} - \boldsymbol{\delta})^{-1}\gamma^{\top}\right)\zeta\\ \iota^{\top}\Xi_{22}\iota &= N\sum_{t=1}^{T}\tilde{\omega}_{t}^{\top}\tilde{\omega}_{t} + T\zeta^{\top}\left(\tilde{X}(NI_{p} - \boldsymbol{\delta})^{-1}\tilde{X}^{\top}\right)\zeta - \frac{N}{T}\left(\sum_{t=1}^{T}\tilde{\omega}_{t}^{\top}\right)\left(\sum_{t=1}^{T}\tilde{\omega}_{t}\right) \end{split}$$

and

$$\begin{bmatrix} \phi^\top & \iota^\top \end{bmatrix} \begin{bmatrix} \Xi_{11} & \Xi_{12} \\ \Xi_{21} & \Xi_{22} \end{bmatrix} \begin{bmatrix} \phi \\ \iota \end{bmatrix} = N \sum_{t=1}^T \tilde{\omega}_t^\top \tilde{\omega}_t - \frac{N}{T} \left(\sum_{t=1}^T \tilde{\omega}_t^\top \right) \left(\sum_{t=1}^T \tilde{\omega}_t \right) + T \zeta^\top \left(\Omega + \left(\gamma - \tilde{X} \right) (NI_p - \boldsymbol{\delta})^{-1} \left(\gamma - \tilde{X} \right)^\top \right) \zeta.$$

with

$$\Omega + (\gamma - \tilde{X}) (NI_p - \delta)^{-1} (\gamma - \tilde{X})^{\top}$$

$$= \Omega + \begin{bmatrix} 0 & \vec{0} \\ \vec{0}^{\top} & I_p - U_{(2)} (I_k + U_{(2)}^{\top} U_{(2)})^{-1} U_{(2)}^{\top} \end{bmatrix}$$

$$= I_N - U(I_k + U^{\top} U)^{-1} U^{\top},$$

following $U^{\top}U = U_{(1)}^{\top}U_{(1)} + U_{(2)}^{\top}U_{(2)}$ and

$$\begin{split} &U_{(2)}(I_k + U_{(2)}^\top U_{(2)})^{-1} U_{(1)}^\top U_{(1)} (I_k + U^\top U)^{-1} U_{(2)}^\top - U_{(2)} (I_k + U_{(2)}^\top U_{(2)})^{-1} U_{(2)}^\top \\ &= U_{(2)} (I_k + U_{(2)}^\top U_{(2)})^{-1} (U_{(1)}^\top U_{(1)} - I_k - U^\top U) (I_k + U^\top U)^{-1} U_{(2)}^\top \\ &= -U_{(2)} (I_k + U^\top U)^{-1} U_{(2)}^\top. \end{split}$$

Thus,

$$\begin{bmatrix} \phi^\top \ \iota^\top \end{bmatrix} \begin{bmatrix} \Xi_{11} \ \Xi_{12} \\ \Xi_{21} \ \Xi_{22} \end{bmatrix} \begin{bmatrix} \phi \\ \iota \end{bmatrix} = N \sum_{t=1}^T \tilde{\omega}_t^\top \tilde{\omega}_t - \frac{N}{T} \left(\sum_{t=1}^T \tilde{\omega}_t^\top \right) \left(\sum_{t=1}^T \tilde{\omega}_t \right) + T \zeta^\top \left(I_N - U (I_k + U^\top U)^{-1} U^\top \right) \zeta.$$

6. Calculate $\vec{z}^{\top} (\Sigma_e^{-1} - \Sigma_e^{-1} \Gamma (\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1}) \vec{z}$.

$$\begin{split} & \vec{z}^{\intercal} \boldsymbol{\Sigma}_{e}^{-1} \vec{z} - \vec{z}^{\intercal} \boldsymbol{\Sigma}_{e}^{-1} \boldsymbol{\Gamma} (\boldsymbol{\Gamma}^{\intercal} \boldsymbol{\Sigma}_{e}^{-1} \boldsymbol{\Gamma})^{-1} \boldsymbol{\Gamma}^{\intercal} \boldsymbol{\Sigma}_{e}^{-1} \vec{z} \\ &= NT - \sum_{t=1}^{T} \vec{z}_{t}^{\intercal} \boldsymbol{U} (\boldsymbol{I}_{k} + \boldsymbol{U}^{\intercal} \boldsymbol{U})^{-1} \boldsymbol{U}^{\intercal} \vec{z}_{t} - N \sum_{t=1}^{T} \tilde{\boldsymbol{\omega}}_{t}^{\intercal} \tilde{\boldsymbol{\omega}}_{t} + \frac{N}{T} \left(\sum_{t=1}^{T} \tilde{\boldsymbol{\omega}}_{t}^{\intercal} \right) \left(\sum_{t=1}^{T} \tilde{\boldsymbol{\omega}}_{t} \right) - T \boldsymbol{\zeta}^{\intercal} \left(\boldsymbol{I}_{N} - \boldsymbol{U} (\boldsymbol{I}_{k} + \boldsymbol{U}^{\intercal} \boldsymbol{U})^{-1} \boldsymbol{U}^{\intercal} \right) \boldsymbol{\zeta} \end{split}$$

Then maximizing $\vec{z}^{\top} \Sigma_e^{-1} \vec{z} - \vec{z}^{\top} \Sigma_e^{-1} \Gamma(\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1} \vec{z}$ is equivalent to minimizing

following $T\sum_{i=1}^{N}\zeta_{i}^{2}=\frac{2N}{T}\sum_{t=1}^{T}\left(T+1-2t\right)\tilde{\omega}_{t,1}$ from the proof of Theorem 1.

Next we show the second step. It is the same as the Sketch of Proof. We separate minimize (a), (b) and (c) and get the sufficient conditions. For (c), note that

$$\begin{split} &\sum_{t=1}^T z_t^\top U(I_k + U^\top U)^{-1} U^\top z_t - \frac{1}{T} \left(\sum_{t=1}^T z_t^\top\right) U(I_k + U^\top U)^{-1} U^\top \left(\sum_{t=1}^T z_t\right) \\ &= \vec{z}^\top \left(\left(I_T - \frac{1}{T} \vec{1} \vec{1}^\top\right) \otimes U(I_k + U^\top U)^{-1} U^\top\right) \vec{z} \coloneqq \vec{z}^\top M_U \vec{z} \end{split}$$

 $I_T - \frac{1}{T} \vec{1} \vec{1}^{\top}$ is a positive semi-definite matrix with one eigenvalue to be 0 and the corresponding eigenvector to be $\vec{1}$. Thus, $\vec{z}^{\top} M_U \vec{z} \geq 0$ for all \vec{z} and the minimum value is attained when $U^{\top} z_t$ is the same for all t, which is equivalent to $Z^{\top} U = \vec{1} \mu_U^{\top} \in \mathbb{R}^{T \times k}$ for some $\mu_U \in \mathbb{R}^k$. For (a), from Theorem 1, the optimal solution of $T \sum_{i=1}^{N} \zeta_i^2 + N \sum_{t=1}^{T} \omega_t^2 - \frac{N}{T} \left(\sum_{t=1}^{T} \omega_t \right)^2$ (equals $(\vec{\omega}^{(1)})^{\top} P_{\vec{1}} \vec{\omega}^{(1)} + 2\vec{d}^{\top} \vec{\omega}^{(1)}$) is $\frac{1}{N} \sum_{i=1}^{N} z_{it} = \frac{2t-1-T}{T}$. For (b), it depends on the cross-sectional average of z_{it} weighted by X_i and is minimized when Z satisfies $\frac{1}{N} \sum_{i=1}^{N} X_i z_{it} = \mu_X$, equivalently $\vec{\omega}^{(j+1)} = \mu_{X,j} \cdot \vec{1}$, where $\mu_{X,j}$ is the j-th entry in μ_X . Hence, we get the sufficient conditions for the optimal solution

$$\frac{1}{N} \sum_{i=1}^{N} z_{it} = \frac{2t - 1 - T}{T}, \qquad \frac{1}{N} \sum_{i=1}^{N} X_i z_{it} = \mu_X, \quad \frac{1}{N} \sum_{i=1}^{N} u_i z_{it} = \mu_U, \text{ for all } t$$

for some $\mu_X \in \mathbb{R}^p$ and $\mu_U \in \mathbb{R}^k$.

Last is to show the third step. Since there are G strata for (X_i, u_i) , that is, G different values for (X_i, u_i) . Then $\omega_t = \sum_{g=1}^G p_g \omega_{g,t}$ and $\omega_t = \sum_{g=1}^G p_g x_{gj} \omega_{g,t}$ for $j = 1, \dots, r$. We have

(b) =
$$N \sum_{t=1}^{T} \sum_{j=1}^{r} \left(\sum_{g=1}^{G} p_g x_{gj} \omega_{g,t} \right)^2 - \frac{N}{T} \sum_{j=1}^{r} \left(\sum_{t=1}^{T} \sum_{g=1}^{G} p_g x_{gj} \omega_{g,t} \right)^2$$

When $\omega_{g,t} = \frac{2t-1-T}{T}$ for all t and g, given X is orthogonal to $\vec{1}$, we have $\sum_{g=1}^{G} p_g x_{gj} \omega_{g,t} = 0$ and then (b) is minimized with value zero. Similarly, we can show that when $\omega_{g,t} = \frac{2t-1-T}{T}$ for all t and g, given U is orthogonal to $\vec{1}$, (c) is minimized with value zero. Moreover, $\omega_t = \sum_{g=1}^{G} p_g \omega_{g,t} = \frac{2t-1-T}{T}$ so (a) = 0 is minimized as well. Hence, the sum of (a), (b) and (c) is minimized.

Proof of Theorem 4 Denote

$$f(\omega_t, \omega_{g,t}) = \sum_{t=1}^T \omega_t^2 - \frac{1}{T} \left(\sum_{t=1}^T \omega_t \right)^2 + 2 \sum_{t=1}^T \frac{T+1-2t}{T} \omega_t$$
$$+ \sum_{t=1}^T \sum_{j=1}^r \left(\sum_{g=1}^G p_g x_{gj} \omega_{g,t} \right)^2 - \frac{1}{T} \sum_{j=1}^r \left(\sum_{t=1}^T \sum_{g=1}^G p_g x_{gj} \omega_{g,t} \right)^2$$

The variance of τ is denoted as $g(\omega_t, \omega_{g,t})$ and

$$g(\omega_t, \omega_{g,t}) \coloneqq \operatorname{Var}(\hat{\tau}) = \frac{\sigma^2}{\vec{z}^\top (I - \Gamma(\Gamma^\top \Gamma)^{-1} \Gamma^\top) \vec{z}} = \frac{\sigma^2}{N(T - f(\omega_t, \omega_{g,t}))}$$

We evaluate $f(\omega_t, \omega_{g,t})$ at the optimal solution in Theorem 3 and has

$$f(\omega_t^*, \omega_{g,t}^*) = -\sum_{t=1}^T \left(\frac{T+1-2t}{T}\right)^2 = -\frac{(T+1)(T-1)}{3T},$$

where $\omega_t^* = \omega_{g,t}^* = \frac{2t-1-T}{T}$. When we use the nearest rounding rule to get a feasible Z^{rnd} , the corresponding $\omega_t, \omega_{g,t}$ are denoted as $\omega_t^{\text{rnd}}, \omega_{g,t}^{\text{rnd}}$. We have $|\omega_{g,t}^{\text{rnd}} - \omega_{g,t}^*| \leq \frac{1}{|\mathcal{O}_g|}$, $|\sum_{t=1}^T \omega_{g,t}^{\text{rnd}}| \leq \frac{1}{|\mathcal{O}_g|}$,

$$|\tilde{w}_t - w_t^*| = |\sum_{g=1}^G p_g (\omega_{g,t}^{\text{rnd}} - \omega_{g,t}^*)| \leq \sum_{g=1}^G p_g |\omega_{g,t}^{\text{rnd}} - \omega_{g,t}^*| \leq \sum_{g=1}^G \frac{p_g}{N_{\min}} = \frac{1}{N_{\min}}$$

and $\left|\sum_{t=1}^{T} \omega_t^{\text{rnd}}\right| \leq \frac{1}{N_{\min}}$. We have

$$\sum_{t=1}^{T} \left(\omega_t^{\text{rnd}} + \frac{T+1-2t}{T} \right)^2 - \frac{1}{T} \left(\sum_{t=1}^{T} \omega_t^{\text{rnd}} \right)^2 - \sum_{t=1}^{T} \left(\frac{T+1-2t}{T} \right)^2 \leq \frac{T}{N_{\min}^2} - \frac{(T+1)(T-1)}{3T}$$

from Theorem 2.

$$\begin{split} & |\sum_{g=1}^{G} p_g x_{g,j} w_{g,t}| \\ & = \Big| \sum_{g=1}^{G} p_g x_{gj} \left(w_{g,t} + \frac{T+1-2t}{T} \right) \Big| = \Big| \sum_{g: x_{gj} > 0} p_g x_{gj} \left(w_{g,t} + \frac{T+1-2t}{T} \right) + \sum_{g: x_{gj} < 0} p_g x_{gj} \left(w_{g,t} + \frac{T+1-2t}{T} \right) \Big| \\ & \leq \max \left(\Big| \sum_{g: x_{gj} > 0} p_g x_{gj} \left(w_{g,t} + \frac{T+1-2t}{T} \right) \Big|, \Big| \sum_{g: x_{gj} < 0} p_g x_{gj} \left(w_{g,t} + \frac{T+1-2t}{T} \right) \Big| \right) \leq \frac{x_{j,\text{max}}}{N_{\text{min}}}, \end{split}$$

and

$$\sum_{t=1}^{T} \sum_{j=1}^{r} \left(\sum_{g=1}^{G} p_g x_{gj} \omega_{g,t} \right)^2 - \frac{1}{T} \sum_{j=1}^{r} \left(\sum_{t=1}^{T} \sum_{g=1}^{G} p_g x_{gj} \omega_{g,t} \right)^2 \le \frac{T \sum_{j=1}^{r} x_{j,\max}^2}{N_{\min}^2}$$

We have

$$\frac{g(\omega_t^{\text{rnd}},\omega_{g,t}^{\text{rnd}})}{g(\omega_t^*,\omega_{g,t}^*)} = \frac{T - f(\omega_t^*,\omega_{g,t}^*)}{T - f(\omega_t^{\text{rnd}},\omega_{g,t}^{\text{rnd}})} \leq \frac{(4T^2 - 1)/(3T)}{(4T^2 - 1)/(3T) - T(1 + \sum_{j=1}^r x_{j,\max}^2)/N_{\min}^2} \leq \frac{1}{1 - (1 + \sum_{j=1}^r x_{j,\max}^2)/N_{\min}^2} \leq \frac{1}{$$

following $3T^2 \leq 4T^2 - 1$ for all T. Note that $f(\omega_t^*, \omega_{g,t}^*)$ is the optimal solution of the convex relaxation problem so $f(\omega_t^*, \omega_{g,t}^*)$ is the lower bound for the minimum objective function value of the original integer programming problem. Thus, $g(\omega_t^*, \omega_{g,t}^*) \leq \operatorname{Var}_{Z^*}(\hat{\tau})$ and we have

$$\operatorname{Var}_{Z^{\operatorname{rnd}}}(\hat{\tau}) \leq \frac{1}{1 - (1 + \sum_{j=1}^{r} x_{j,\max}^2) / N_{\min}^2} \operatorname{Var}_{Z^*}(\hat{\tau})$$

Appendix E: Extension to Carryover Effects

Proof of Lemmas 4 and 5 Let $\vec{z}^{(j)} := \begin{bmatrix} \vec{z}_j^\top \vec{z}_{j+1}^\top \cdots \vec{z}_{T-\ell-1+j}^\top \end{bmatrix}^\top$ for $j = 1, \dots, \ell+1$. Then $\mathcal{Z} = 1$ $[\vec{z}^{(1)}\ \vec{z}^{(2)}\ \cdots\ \vec{z}^{(\ell+1)}]$. Furthermore, we decompose $\mathrm{Var}(\hat{\vec{\tau}})$

$$\Theta = \mathcal{Z}^{\top} (I_{N(T-\ell)} - \Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top}) \mathcal{Z} = \underbrace{\mathcal{Z}^{\top}\mathcal{Z}}_{\Theta_1} - \underbrace{\mathcal{Z}^{\top}\Gamma(\Gamma^{\top}\Gamma)^{-1}\Gamma^{\top}\mathcal{Z}}_{\Theta_2}$$

and simplify Θ_1 and Θ_2 separately. Let us first analyze $\Theta_1 = \mathcal{Z}^{\top} \mathcal{Z}$. For the (j,j)-th entry in Θ_1 , we have $\Theta_{1,jj} = N(T - \ell)$. For the (j,m)-th entry in $\mathcal{Z}^{\top}\mathcal{Z}$ (j < m), we have

$$\Theta_{1,jm} = N \left[(T - \ell) + \sum_{t=j}^{m-1} (\omega_t - \omega_{T-\ell+t}) \right].$$

Next, let use consider
$$\Theta_2 = \mathcal{Z}^{\top} \Gamma(\Gamma^{\top} \Gamma)^{-1} \Gamma^{\top} \mathcal{Z}$$
. We have
$$(\Gamma^{\top} \Gamma)^{-1} = \begin{bmatrix} \Xi_{11} \ \Xi_{12} \\ \Xi_{21} \ \Xi_{22} \end{bmatrix} = \begin{bmatrix} M & -\frac{1}{N} M M_{\vec{1}} \\ -\frac{1}{N} M_{\vec{1}}^{\top} M \ \frac{1}{N} I_{T-\ell} + \frac{1}{N^2} M_{\vec{1}}^{\top} M M_{\vec{1}} \end{bmatrix} \in \mathbb{R}^{(N+T-\ell-1)\times(N+T-\ell-1)},$$

where $\Xi_{11} = M$, $\Xi_{12} = -\frac{1}{N}MM_{\vec{1}}$, $\Xi_{21} = -\frac{1}{N}M_{\vec{1}}^{\top}M$, $\Xi_{22} = \frac{1}{N}I_{T-\ell} + \frac{1}{N^2}M_{\vec{1}}^{\top}MM_{\vec{1}}$ and

$$M = \frac{1}{T - \ell} \left(I_{N-1} - \frac{1}{N} \vec{1} \vec{1}^{\top} \right)^{-1} = \frac{1}{T - \ell} (I_{N-1} + \vec{1} \vec{1}^{\top}) \in \mathbb{R}^{(N-1) \times (N-1)}, \quad M_{\vec{1}} = \left[\vec{1} \cdots \vec{1} \right] \in \mathbb{R}^{(N-1) \times (T-\ell)}.$$

Furthermore

$$(\vec{z}^{(j)})^{\top} \Gamma = \left[\sum_{t=j}^{T-\ell-1+j} z_{1t} \cdots \sum_{t=j}^{T-\ell-1+j} z_{N-1,t} \sum_{i=1}^{N} z_{ij} \cdots \sum_{i=1}^{N} z_{i,T-\ell-1+j} \right]$$

$$= \left[(T-\ell)\zeta_{1}^{(j)} \cdots (T-\ell)\zeta_{N-1}^{(j)} N\omega_{j} \cdots N\omega_{T-\ell-1+j} \right],$$

where $\zeta_i^{(j)} = \frac{1}{T-\ell} \sum_{t=j}^{T-\ell-1+j} z_{1t}$ for $i = 1, 2, \dots, N$ and $j = 1, 2, \dots, \ell+1$, and $\omega_t = \frac{1}{N} \sum_{i=1}^{T} z_{it}$ for $t = 1, 2, \dots, T$. We let $\alpha^{(j)} = \left[(T-\ell)\zeta_1^{(j)} \cdots (T-\ell)\zeta_{N-1}^{(j)} \right]^{\top}$ and $\beta^{(j)} = \left[Ny_j \cdots Ny_{T-\ell-1+j} \right]^{\top}$. Then

$$\mathcal{Z}^{\top}\Gamma = \begin{bmatrix} (\alpha^{(1)})^{\top} & (\beta^{(1)})^{\top} \\ \vdots & \vdots \\ (\alpha^{(\ell+1)})^{\top} & (\beta^{(\ell+1)})^{\top} \end{bmatrix}$$

and the (j,m)-th term of $\Theta_2 = \mathcal{Z}^\top \Gamma(\Gamma^\top \Gamma)^{-1} \Gamma^\top \mathcal{Z}$ ha

$$\Theta_{2,jm} = (\alpha^{(j)})^{\top} \Xi_{11} \alpha^{(m)} + (\alpha^{(j)})^{\top} \Xi_{12} \beta^{(m)} + (\beta^{(j)})^{\top} \Xi_{21} \alpha^{(m)} + (\beta^{(j)})^{\top} \Xi_{22} \beta^{(m)} \in \mathbb{R}^{(\ell+1) \times (N+T-\ell-1)}.$$

For each term in $\Theta_{2,im}$,

$$\begin{split} (\alpha^{(j)})^{\top} \Xi_{11} \alpha^{(m)} &= \frac{1}{T - \ell} (\alpha^{(j)})^{\top} (I_{N-1} + \vec{1}\vec{1}^{\top}) \alpha^{(m)} \\ &= (T - \ell) \sum_{i=1}^{N-1} \zeta_{i}^{(j)} \zeta_{i}^{(m)} + (T - \ell) \left(\sum_{i=1}^{N-1} \zeta_{i}^{(j)} \right) \left(\sum_{i=1}^{N-1} \zeta_{i}^{(m)} \right) \\ (\alpha^{(j)})^{\top} \Xi_{12} \beta^{(m)} &= -\frac{1}{T - \ell} (\alpha^{(j)})^{\top} (I_{N-1} + \vec{1}\vec{1}^{\top}) \vec{1} \left(\sum_{t=m}^{T - \ell - 1 + m} \omega_{t} \right) = -\frac{N}{T - \ell} (\alpha^{(j)})^{\top} \vec{1} \left(\sum_{t=m}^{T - \ell - 1 + m} \omega_{t} \right) \\ &= -N \left(\sum_{i=1}^{N-1} \zeta_{i}^{(j)} \right) \left(\sum_{t=m}^{T - \ell - 1 + m} \omega_{t} \right) \\ (\beta^{(j)})^{\top} \Xi_{21} \alpha^{(m)} &= -\frac{1}{T - \ell} (\alpha^{(m)})^{\top} (I_{N-1} + \vec{1}\vec{1}^{\top}) \vec{1} \left(\sum_{t=j}^{T - \ell - 1 + j} \omega_{t} \right) = -\frac{N}{T - \ell} (\alpha^{(m)})^{\top} \vec{1} \left(\sum_{t=j}^{T - \ell - 1 + j} \omega_{t} \right) \\ &= -N \left(\sum_{i=1}^{N-1} \zeta_{i}^{(m)} \right) \left(\sum_{t=j}^{T - \ell - 1 + j} \omega_{t} \right) \\ (\beta^{(j)})^{\top} \Xi_{22} \beta^{(m)} &= N \sum_{t=j}^{T - \ell - 1 + j} \omega_{t} \omega_{t+m-j} + \frac{N(N-1)}{T - \ell} \left(\sum_{t=j}^{T - \ell - 1 + j} \omega_{t} \right) \left(\sum_{t=m}^{T - \ell - 1 + m} \omega_{t} \right) \end{split}$$

From the definition of $\zeta_i^{(j)}$ and ω_t , we have $(T-\ell)\sum_{i=1}^N \zeta_i^{(j)} = N\sum_{t=j}^{T-\ell-1+j} \omega_t$ for $j=1,2,\cdots,\ell+1$ and $\sum_{i=1}^{N-1} \zeta_i^{(j)} = \frac{N}{T-\ell}\sum_{t=j}^{T-\ell-1+j} \omega_t - \zeta_N^{(j)}$. Using this property, we have

$$\Theta_{2,jm} = (T - \ell) \sum_{i=1}^{N} \zeta_i^{(j)} \zeta_i^{(m)} - \frac{N}{T - \ell} \left(\sum_{t=j}^{T-\ell-1+j} \omega_t \right) \left(\sum_{t=m}^{T-\ell-1+m} \omega_t \right) + N \sum_{t=j}^{T-\ell-1+j} \omega_t \omega_{t+m-j}.$$
 (32)

Recall the definition $\zeta_i^{(j)} = \frac{1}{T-\ell} \sum_{t=j}^{T-\ell-1+j} z_{1t}$, there are T+1 different values that $\zeta_i^{(j)} \zeta_i^{(m)}$ can take, denoted as $v_0^{(j,m)}, v_1^{(j,m)}, \cdots, v_T^{(j,m)}$, where $v_t^{(j,m)}$ denotes the value of $\zeta_i^{(j)} \zeta_i^{(m)}$ when unit i starts to get the treatment at time period T+1-t (and $v_0^{(j,m)}$ represents the value of $\zeta_i^{(j)} \zeta_i^{(m)}$ when unit i stays in the control group for all time periods). Without loss of generality, we assume $j \leq m$ and have

$$v_t^{(j,m)} = \begin{cases} 1 & t \le \ell + 1 - m \\ -\left(-1 + \frac{2(t-1-\ell+m)}{T-\ell}\right) & \ell + 1 - m < t \le \ell + 1 - j \\ \left(-1 + \frac{2(t-1-\ell+m)}{T-\ell}\right) \left(-1 + \frac{2(t-1-\ell+j)}{T-\ell}\right) & \ell + 1 - j < t \le T + 1 - k \\ \left(-1 + \frac{2(t-1-\ell+j)}{T-\ell}\right) & T + 1 - m < t \le T + 1 - j \\ 1 & T + 1 - j < t \end{cases}$$

Given ω_t , there are $\frac{N(1+\omega_1)}{2}$, $\frac{N(1+\omega_2)}{2}$, \cdots $\frac{N(1+\omega_T)}{2}$ treated units in time period $1, 2, \cdots, T$. It is equivalent to having $\frac{N(1+\omega_1)}{2}$, $\frac{N(\omega_2-\omega_1)}{2}$, \cdots , $\frac{N(\omega_T-\omega_{T-1})}{2}$ untreated units to start the treatment in time period $1, 2, \cdots, T$ and leaving $\frac{N(1-\omega_T)}{2}$ units in the control group in the end.

$$\begin{split} \sum_{i=1}^{N} \zeta_{i}^{(j)} \zeta_{i}^{(m)} &= N \left[\frac{1+\omega_{1}}{2} \cdot v_{T}^{(j,m)} + \frac{\omega_{2}-\omega_{1}}{2} v_{T-1}^{(j,m)} + \dots + \frac{1-\omega_{T}}{2} \cdot v_{0}^{(j,m)} \right] \\ &= N \left[1 + \frac{v_{T}^{(j,m)} - v_{T-1}^{(j,m)}}{2} \omega_{1} + \frac{v_{T-1}^{(j,m)} - v_{T-2}^{(j,m)}}{2} \omega_{2} + \dots + \frac{v_{1}^{(j,m)} - v_{0}^{(j,m)}}{2} \omega_{T} \right], \end{split}$$

following $v_0^{(j,m)} = v_T^{(j,m)} = 1$. Then we can rewrite $\Theta_{2,jm}$ in terms of ω_t and $v_t^{(j,m)}$

$$\Theta_{2,jm} = N(T-\ell) + \frac{N(T-\ell)}{2} \sum_{t=1}^{T} (\upsilon_{T+1-t}^{(j,m)} - \upsilon_{T-t}^{(j,m)}) \omega_t - \frac{N}{T-\ell} \left(\sum_{t=j}^{T-\ell-1+j} \omega_t \right) \left(\sum_{t=m}^{T-\ell-1+m} \omega_t \right) + N \sum_{t=j}^{T-\ell-1+j} \omega_t \omega_t \omega_{t+m-j}.$$

Recall $\Theta = \Theta_1 - \Theta_2$, for the diagonal entries in Θ ,

$$\Theta_{jj} = -N \left[\sum_{t=j}^{T-\ell-1+j} \omega_t^2 - \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell-1+j} \omega_t \right)^2 + \frac{T-\ell}{2} \sum_{t=1}^T (\upsilon_{T+1-t}^{(j,j)} - \upsilon_{T-t}^{(j,j)}) \omega_t \right]$$

and for the off-diagonal entries in Θ and for j < m,

$$\Theta_{jm} = -N \left[\sum_{t=j}^{T-\ell-1+j} \omega_t \omega_{t+m-j} - \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell-1+j} \omega_t \right) \left(\sum_{t=m}^{T-\ell-1+m} \omega_t \right) + \frac{T-\ell}{2} \sum_{t=1}^{T} (\upsilon_{T+1-t}^{(j,m)} - \upsilon_{T-t}^{(j,m)}) \omega_t - \sum_{t=j}^{m-1} (\omega_t - \omega_{T-\ell+t}) \right]$$

and for j > m,

$$\Theta_{jm} = \Theta_{mj}$$
.

Thus, for the **T**-optimal design, the optimization problem is

$$\min_{\omega_{t}} \sum_{j=1}^{\ell+1} \left[\sum_{t=j}^{T-\ell-1+j} \omega_{t}^{2} - \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell-1+j} \omega_{t} \right)^{2} + \frac{T-\ell}{2} \sum_{t=1}^{T} (v_{T+1-t}^{(j,j)} - v_{T-t}^{(j,j)}) \omega_{t} \right]
\text{s.t.} \quad -1 \le \omega_{t} \le 1, \quad t = 1, 2, \dots, T
\omega_{t} \le \omega_{t+1} \quad t = 1, 2, \dots, T - 1$$
(33)

and we can show $\frac{T-\ell}{2} = \frac{2(T-\ell-1+2j-2t)}{T-\ell}$

For the **D**-optimal treatment design, the optimization problem is

$$\min_{\omega_t} -det(\Theta)
\text{s.t.} \quad -1 \le \omega_t \le 1, \quad t = 1, 2, \dots, T
\omega_t \le \omega_{t+1}, \quad t = 1, 2, \dots, T - 1$$
(34)

The following lemma provides an optimal solution for the quadratic program (25) in Lemma 4. This lemma provides an important intermediate step to show Theorem 5.

LEMMA 6 (Carryover Effects without Covariates). Suppose Assumptions 1-2 hold, the potential outcome follows model (14) and does not have observed or latent covariates, i.e., $Y_{it}(z_{it}) = \alpha_i + \beta_t + \tau_1 z_{it} + \cdots + \tau_{\ell+1} z_{i,t-\ell} + \varepsilon_{it}$, and $T > \frac{\ell^3 + 13\ell^2 + 7\ell + 3}{8\ell}$, then an optimal solution for the quadratic program (25) is $\vec{\omega}^* \in \mathbb{R}^T$ that is defined in Eq. (17).

Proof of Lemma 6 Denote the objective in the quadratic program (25) in Lemma 4 as $f(\vec{\omega}) \coloneqq \sum_{j=1}^{\ell+1} \left[\sum_{t=j}^{T-\ell-1+j} \omega_t^2 - \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell-1+j} \omega_t \right)^2 + \sum_{t=1}^{T-\ell-1+j} \frac{2(T-\ell-1+2j-2t)}{T-\ell} \omega_t \right].$ The Lagrangian of $f(\vec{\omega})$ is

$$\mathcal{L}(\vec{\omega}, \vec{\lambda}, \vec{\kappa}, \vec{\iota}) = \sum_{j=1}^{\ell+1} \left[\sum_{t=j}^{T-\ell-1+j} \omega_t^2 - \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell-1+j} \omega_t \right)^2 + \sum_{t=1}^{T} \frac{2(T-\ell-1+2j-2t)}{T-\ell} \omega_t \right] + \sum_{t=1}^{T} \lambda_t (-1-\omega_t) + \sum_{t=1}^{T} \kappa_t (\omega_t - 1) + \sum_{t=1}^{T-1} \iota_t (\omega_t - \omega_{t+1})$$

The KKT conditions of $\mathcal{L}(\vec{\omega}, \vec{\lambda}, \vec{\kappa}, \vec{\iota})$ are

$$\frac{\partial \mathcal{L}}{\partial \omega_t} = t\omega_t - \frac{\sum_{j=1}^t s_j}{T - \ell} + \frac{(T - \ell - t)t}{T - \ell} - \lambda_t + \kappa_t + \iota_t - \iota_{t-1} = 0, \quad t \le \ell$$
(35)

$$\frac{\partial \mathcal{L}}{\partial \omega_t} = (\ell + 1)\omega_t - \frac{\sum_{j=1}^{\ell+1} s_j}{T - \ell} + \frac{(\ell + 1)(T + 1 - 2t)}{T - \ell} - \lambda_t + \kappa_t + \iota_t - \iota_{t-1} = 0, \quad \ell < t \le T - \ell$$
 (36)

$$\frac{\partial \mathcal{L}}{\partial \omega_{t}} = (T+1-t)\omega_{t} - \frac{\sum_{j=1}^{T+1-t} s_{j}}{T-\ell} + \frac{(T-\ell-1)(T+1-t)}{T-\ell} - \lambda_{t} + \kappa_{t} + \iota_{t} - \iota_{t-1} = 0, \ t > T-\ell
\lambda_{t}(-1-\omega_{t}) = 0, \quad \kappa_{t}(\omega_{t}-1) = 0, \quad \iota_{t}(\omega_{t}-\omega_{t+1}) = 0
-1 \le \omega_{t} \le 1, \quad \omega_{t} \le \omega_{t+1}, \quad \lambda_{t} \ge 0, \quad \kappa_{t} \ge 0, \quad \iota_{t} \ge 0$$
(37)

where $s_j = \sum_{t=j}^{T-\ell-1+j} \omega_t$ for $j = 1, \dots, \ell+1$ and $\iota_0 = 0$.

The Hessian of $f(\vec{\omega})$ is positive semi-definite. Any solution that satisfies the KKT conditions is optimal.

First we can show the optimal solution is symmetric with respect to the origin. The proof is as follows. If $\vec{\omega}^{\ddagger}$ is the optimal solution of the quadratic program (25), then we can show $\vec{\omega}^{\dagger} = \left[-\omega_T^{\ddagger} - \omega_{T-1}^{\ddagger} \cdots - \omega_1^{\ddagger}\right]$ has the same value in the objective function as $\vec{\omega}^{\ddagger}$ because $\sum_{j=1}^{\ell+1} \sum_{t=1}^{T-\ell-1+j} \frac{2(T-\ell-1+2j-2t)}{T-\ell} \omega_t$ in $f(\vec{\omega})$ is symmetric with respect to the origin and similarly for the other two terms in $f(\vec{\omega})$. Since quadratic program (25) is convex, $f\left(\frac{\vec{\omega}^{\ddagger}+\vec{\omega}^{\dagger}}{2}\right) \leq \frac{1}{2} \left[f(\vec{\omega}^{\ddagger}) + f(\vec{\omega}^{\dagger})\right] = f(\vec{\omega}^{\ddagger})$. Then if $\vec{\omega}^{\ddagger}$ is optimal, $\vec{\omega}^{\dagger} = \vec{\omega}^{\ddagger}$.

Now we can focus on the ω that satisfies $\left[\omega_1 \ \omega_2 \ \cdots \ \omega_T\right] = \left[-\omega_T \ -\omega_{T-1} \ \cdots \ -\omega_1\right]$. From the definition of $s_j = \sum_{t=j}^{T-\ell-1+j} \omega_t$, we have $s_j = -s_{\ell+1-j}$. If ℓ is even, $s_{\ell/2+1} = 0$.

Now we are going to verify $\vec{\omega}^* = \left[\omega_1^* \ \omega_2^* \cdots \omega_T^*\right]$ defined in Eq. (17) satisfies the KKT conditions with feasible $\vec{\lambda}, \vec{\kappa}, \vec{\iota}$.

- 1. ω_t^* for $\ell < t \le T \ell$: $\omega_t^* = -1 + \frac{2t (\ell + 1)}{T \ell}$ satisfies Eq. (36) with $\lambda_t = \kappa_t = \iota_t = 0$ and $\iota_\ell = 0$.
- 2. ω_t^* for $t \leq \ell$: Given $\omega_t = -\omega_{T+1-t}$. We can simplify s_j to

$$s_j = \begin{cases} \sum_{\substack{t=j \\ t=T-j}}^{\ell+1-j} \omega_t & \text{ for } j=1,\cdots,\lfloor (\ell+1)/2 \rfloor \\ \sum_{\substack{t=T-j \\ t=T-j}}^{T-\ell+j} \omega_t & \text{ for } j=\lfloor (\ell+1)/2 \rfloor +1,\cdots,\ell+1 \end{cases}.$$

As an example, when $\ell=2$, we have $s_1=\omega_1+\omega_2$, $s_2=0$ and $s_3=\omega_{T-1}+\omega_T$; when $\ell=3$, we have $s_1=\omega_1+\omega_2+\omega_3$, $s_2=\omega_2$, $s_3=\omega_{T-1}$ and $s_4=\omega_{T-2}+\omega_{T-1}+\omega_T$. Furthermore, $s_j+s_{\ell+2-j}=0$ for $1\leq j\leq \ell+1$. Using this property, for $\lfloor \ell/2 \rfloor < t \leq \ell$, we have $\sum_{j=1}^t s_j = \sum_{j=1}^{\ell+1-t} s_j$.

Next we show when $\omega_t = -1$ for $t \leq \lfloor \ell/2 \rfloor$, there exist some ω_t for $\lfloor \ell/2 \rfloor < t \leq \ell$ and some feasible $\lambda_t, \kappa_t, \iota_t$ that satisfy Eq. (35).

When $\omega_t = -1$ for $t \leq \lfloor \ell/2 \rfloor$, then for $\lfloor \ell/2 \rfloor < t \leq \ell$, $\sum_{j=1}^t s_j = \sum_{j=1}^{\ell+1-t} s_j = \left[\sum_{j=1}^{\ell+1-t} (\lfloor \ell/2 \rfloor + 1 - j) \right] + \min(\ell+1-t,\ell-\lfloor \ell/2 \rfloor) \omega_{\lfloor \ell/2 \rfloor+1} + \cdots + \min(\ell+1-t,2) \omega_{\ell-1} + \min(\ell+1-t,1) \omega_{\ell}$. As an example, when $\ell=2$, $s_2 = -1 + \omega_2$; when $\ell=3$, $s_1 + s_2 = -1 + 2\omega_2 + \omega_3$, $s_1 + s_2 + s_3 = -1 + \omega_2 + \omega_3$. We can rewrite Eq. (35) for $\lfloor \ell/2 \rfloor < t \leq \ell$ in a vectorized form as the following (we will consider Eq. (35) for $t \leq \lfloor \ell/2 \rfloor$ in the later part of this proof)

$$\begin{bmatrix}
\frac{\partial \mathcal{L}}{\partial \omega_{\lfloor \ell/2 \rfloor + 1}} \\
\vdots \\
\frac{\partial \mathcal{L}}{\partial \omega_{\ell}}
\end{bmatrix} = A^{(\ell)} \begin{bmatrix} \omega_{\lfloor \ell/2 \rfloor + 1} \\
\vdots \\
\omega_{\ell} \end{bmatrix} - b^{(\ell)} - \begin{bmatrix} \lambda_{\lfloor \ell/2 \rfloor + 1} \\
\vdots \\
\lambda_{\ell} \end{bmatrix} + \begin{bmatrix} \kappa_{\lfloor \ell/2 \rfloor + 1} \\
\vdots \\
\kappa_{\ell} \end{bmatrix} + \begin{bmatrix} \iota_{\lfloor \ell/2 \rfloor + 1} \\
\vdots \\
\iota_{\ell} \end{bmatrix} - \begin{bmatrix} \iota_{\lfloor \ell/2 \rfloor} \\
\vdots \\
\iota_{\ell-1} \end{bmatrix} = 0,$$
(38)

where $A^{(\ell)}$ and $b^{(\ell)}$ are defined in Eq. (22) and Eq. (23). When $\left[\omega_{\lfloor \ell/2 \rfloor + 1} \cdots \omega_{\ell}\right]^{\top} = (A^{(\ell)})^{-1}b^{(\ell)}$, Eq. (38) holds with $\lambda_t = \kappa_t = \iota_t = 0$ for $t = \lfloor \ell/2 \rfloor + 1, \cdots, \ell$ and $\iota_{\lfloor \ell/2 \rfloor} = 0$. The remaining step is to verify the constraints $-1 \le \omega_t \le -1 + \frac{\ell+1}{T-\ell}$ and $\omega_t \le \omega_{t+1}$ hold if $\left[\omega_{\lfloor \ell/2 \rfloor + 1} \cdots \omega_{\ell}\right]^{\top} = (A^{(\ell)})^{-1}b^{(\ell)}$.

(a) The first step is to show $-A^{(\ell)}\vec{1} \leq b^{(\ell)} \leq (-1 + \frac{\ell+1}{T-\ell})A^{(\ell)}\vec{1}$. Note that the diagonal entries in $A^{(\ell)}$ are positive while the off-diagonal entries in $A^{(\ell)}$ are negative, then $A^{(\ell)}_{t',:}\vec{\omega}_{(\lfloor \ell/2 \rfloor+1):L}$ is increasing in ω_t and decreasing in ω_s for $t'=1,\cdots,\ell-\lfloor \ell/2 \rfloor$, $t=t'+\lfloor \ell/2 \rfloor$ and $s\neq t'+\lfloor \ell/2 \rfloor$. If $-A^{(\ell)}\vec{1} \leq b^{(\ell)} \leq (-1 + \frac{\ell+1}{T-\ell})A^{(\ell)}\vec{1}$ hold, then ω_t defined in $\left[\omega_{\lfloor \ell/2 \rfloor+1}\cdots\omega_\ell\right]^\top=(A^{(\ell)})^{-1}b^{(\ell)}$ is between -1 and $-1 + \frac{\ell+1}{T-\ell}$ for $t=\lfloor \ell/2 \rfloor+1,\cdots,\ell$.

First, let us show $-A^{(\ell)}\vec{1} \leq b^{(\ell)}$, which is equivalent to showing every entry in $A^{(\ell)}\vec{1} + b^{(\ell)}$ is non-negative, that is, for $t' = 1, \dots, \ell - \lfloor \ell/2 \rfloor$, $(A^{(\ell)}\vec{1})_{t'} + b^{(\ell)}_{t'} \geq 0$. If ℓ is even, $\sum_{l=1}^{\ell-t} (\ell - \lfloor \ell/2 \rfloor + 1 - l) = \frac{t(\ell+1-t)}{2}$ and $\sum_{l=1}^{\ell-t} (\lfloor \ell/2 \rfloor + 1 - l) = \frac{t(\ell+1-t)}{2}$. Let $t = t' + \lfloor \ell/2 \rfloor$. We have

$$(A^{(\ell)}\vec{1})_{t'} + b_{t'}^{(\ell)} = t - \frac{1}{T - \ell} \frac{t(\ell + 1 - t)}{2} - t + \frac{t^2}{T - \ell} - \frac{1}{T - \ell} \frac{t(\ell + 1 - t)}{2} = \frac{t(2T - \ell - 1)}{T - \ell} \ge 0.$$

If ℓ is odd, $\sum_{j=1}^{\ell-t} (\ell - \lfloor \ell/2 \rfloor + 1 - j) = \frac{(t+1)(\ell+1-t)}{2}$ and $\sum_{j=1}^{\ell-t} (\lfloor \ell/2 \rfloor + 1 - j) = \frac{(t-1)(\ell+1-t)}{2}$. We have

$$(A^{(\ell)}\vec{1})_{t'} + b_{t'}^{(\ell)} = t - \frac{1}{T-\ell} \frac{(t+1)(\ell+1-t)}{2} - t + \frac{t^2}{T-\ell} - \frac{1}{T-\ell} \frac{(t-1)(\ell+1-t)}{2} = \frac{t(2T-\ell-1)}{T-\ell} \geq 0.$$

Second, let us show $b^{(\ell)} \leq (-1 + \frac{\ell+1}{T-\ell})A^{(\ell)}\vec{1}$, which is equivalent to showing every entry in $(1 - \frac{\ell+1}{T-\ell})A^{(\ell)}\vec{1} + b^{(\ell)}$ is non-positive, that is, for $t' = 1, \dots, \ell - \lfloor \ell/2 \rfloor$, $\left(A^{(\ell)}(1 - \frac{\ell+1}{T-\ell})\vec{1}\right)_{t'} + b^{(\ell)}_{t'} \leq 0$. If ℓ is even

$$\Big(A^{(\ell)}(1-\frac{\ell+1}{T-\ell})\vec{1}\Big)_{t'} + b^{(\ell)}_{t'} = \frac{t(2T-\ell-1)}{T-\ell} - \frac{\ell+1}{T-\ell}\Big(t - \frac{t(\ell+1-t)}{2(T-\ell)}\Big) = \frac{t(T-\ell-1)}{T-\ell}\Big(2 - \frac{1}{2}\frac{\ell+1}{T-\ell}\Big) < 0$$

following $t(T-\ell-1) < 0$ and $2-\frac{1}{2}\frac{\ell+1}{T-\ell} > 0$. If ℓ is odd,

$$\left(A^{(\ell)} (1 - \frac{\ell+1}{T-\ell}) \vec{1} \right)_{t'} + b_{t'}^{(\ell)} = \frac{t(2T-\ell-1)}{T-\ell} - \frac{\ell+1}{T-\ell} \left(t - \frac{(t+1)(\ell+1-t)}{2(T-\ell)} \right) = \frac{t(T-\ell-1)}{T-\ell} \left(2 - \frac{t+1}{2t} \frac{\ell+1}{T-\ell} \right) < 0$$
 following $t(T-\ell-1) < 0$ and $2 - \frac{t+1}{2t} \frac{\ell+1}{T-\ell} > 0$.

(b) We can show $-b_{t'}^{(\ell)}/(A^{(\ell)}\vec{1})_{t'}$ is non-decreasing in t' for $t'=1,\cdots,\ell-\lfloor\ell/2\rfloor$. Note that the diagonal entries in $A^{(\ell)}$ are positive while the off-diagonal entries in $A^{(\ell)}$ are negative, then $A_{t',:}^{(\ell)}\vec{\omega}_{(\lfloor\ell/2\rfloor+1):L}$ is increasing in ω_t and decreasing in ω_s for $t'=1,\cdots,\ell-\lfloor\ell/2\rfloor$, $t=t'+\lfloor\ell/2\rfloor$ and $s\neq t'+\lfloor\ell/2\rfloor$. If $-b_{t'}^{(\ell)}/(A^{(\ell)}\vec{1})_{t'}$ is non-decreasing in t', then ω_t is non-decreasing in t, where $t=t'+\lfloor\ell/2\rfloor$.

Let $c_{t'}$ be the $c_{t'}$ that satisfies $(A^{(\ell)}(-1+\frac{c_{t'}}{T-\ell})\vec{1})_{t'}=b_{t'}^{(\ell)}$ and let $t=t'+\lfloor \ell/2 \rfloor$. If ℓ is even, we have

$$\begin{split} \frac{t(2T-\ell-1)}{T-\ell} &= \frac{c_{t'}}{T-\ell} \left(t - \frac{1}{T-\ell} \frac{t(\ell+1-t)}{2} \right) \\ \Leftrightarrow 2T-\ell-1 &= c_{t'} \frac{t+2T-3\ell-1}{2(T-\ell)}. \end{split}$$

Since $\frac{\partial(2T-\ell-1)}{\partial t} = 2$, $\frac{\partial \frac{t+2T-3\ell-1}{2(T-\ell)}}{\partial t} = \frac{1}{2(T-\ell)}$, and $2 > \frac{1}{2(T-\ell)}$, we have $c_{t'}$ increases in t and t'. This implies $-b_{t'}^{(\ell)}/(A^{(\ell)}\vec{1})_{t'}$ is non-decreasing in t' for even ℓ .

If ℓ is odd, we have

$$\begin{split} \frac{t(2T-\ell-1)}{T-\ell} &= \frac{c_{t'}}{T-\ell} \left(t - \frac{1}{T-\ell} \frac{(t+1)(\ell+1-t)}{2} \right) \\ \Leftrightarrow 2T-\ell-1 &= c_{t'} \left(1 + \frac{(t+1)(T-\ell-1)}{2t(T-\ell)} \right). \end{split}$$

Since $\frac{\partial(2T-\ell-1)}{\partial t}=2$, $\frac{\partial\frac{t+2T-3\ell-1}{2(T-\ell)}}{\partial t}\leq\frac{\ell+3}{\ell+1}\frac{1}{T-\ell}$, and $2>\frac{\ell+3}{(T-\ell)(\ell+1)}$, we have $c_{t'}$ increases in t and t'. This again implies $-b_{t'}^{(\ell)}/(A^{(\ell)}\vec{1})_{t'}$ is non-decreasing in t' for odd ℓ .

We have verified that for $\lfloor \ell/2 \rfloor < t \le \ell$, ω_t defined in $\left[\omega_{\lfloor \ell/2 \rfloor + 1} \cdots \omega_\ell\right]^\top = (A^{(\ell)})^{-1}b^{(\ell)}$ satisfies the KKT conditions. The remaining step is to verify for $t \le \lfloor \ell/2 \rfloor$, ω_t defined as $\omega_t = -1$ satisfies the KKT conditions. When $\omega_t = -1$, constraints $-1 \le \omega_t \le 1$, $\omega_t \le \omega_{t+1}$ for $t \le \lfloor \ell/2 \rfloor$ and $\omega_{\lfloor \ell/2 \rfloor} \le \omega_{t+\lfloor \ell/2 \rfloor + 1}$ are satisfied. We only need to verify that we can find feasible $\lambda_t, \kappa_t, \iota_t$ to satisfy Eq. (35). Since $\omega_t = -1$, from complementary slackness, $\kappa_t = 0$. Plug $\omega_t = -1$ into Eq. (35), we have

$$\lambda_{1} - \iota_{1} = -\frac{1+s_{1}}{T-\ell}$$

$$\lambda_{t} - \iota_{t} + \iota_{t-1} = -\frac{t^{2} + \sum_{j=1}^{t} s_{j}}{T-\ell} \text{ for } t = 2, \dots \lfloor \ell/2 \rfloor - 1$$

$$\lambda_{t} + \iota_{t-1} = -\frac{t^{2} + \sum_{j=1}^{t} s_{j}}{T-\ell} \text{ for } t = \lfloor \ell/2 \rfloor$$

We only need to verify $-\frac{t^2 + \sum_{j=1}^t s_j}{T - \ell} \ge 0$ for $t = \lfloor \ell/2 \rfloor$ as for the other conditions, $\lambda_1 - \iota$ and $\lambda_t - \iota_t + \iota_{t-1}$ can take any value by properly choosing λ_t and ι .

Note that $\frac{1}{T-\ell}\sum_{j=1}^{\lfloor\ell/2\rfloor}s_j=\frac{1}{T-\ell}\sum_{j=1}^{L+1-\lfloor\ell/2\rfloor}s_j=(\ell+1-\lfloor\ell/2\rfloor)(\omega_{\ell+1-\lfloor\ell/2\rfloor}+1)-\frac{(\ell+1-\lfloor\ell/2\rfloor)^2}{T-\ell}$. Furthermore, if we can show $\omega_{\ell+1-\lfloor\ell/2\rfloor}+1\leq \frac{\ell+1}{T-\ell}\frac{1}{\ell+1-\lfloor\ell/2\rfloor}$ for even ℓ and $\omega_{\ell+1-\lfloor\ell/2\rfloor}+1\leq \frac{\ell+1}{T-\ell}\frac{2}{\ell+1-\lfloor\ell/2\rfloor}$ for odd ℓ , then we have

$$\begin{aligned} &-\frac{\lfloor \ell/2 \rfloor^2}{T-\ell} - (\ell+1 - \lfloor \ell/2 \rfloor)(\omega_{\ell+1-\lfloor \ell/2 \rfloor} + 1) + \frac{(\ell+1 - \lfloor \ell/2 \rfloor)^2}{T-\ell} \\ &= \frac{(\ell+1 - 2\lfloor \ell/2 \rfloor)(\ell+1)}{T-\ell} - (\ell+1 - \lfloor \ell/2 \rfloor)(\omega_{\ell+1-\lfloor \ell/2 \rfloor} + 1) \ge 0 \end{aligned}$$

and therefore $-\frac{t^2 + \sum_{j=1}^t s_j}{T - \ell} \ge 0$.

Next is to show " $\omega_{\ell+1-\lfloor\ell/2\rfloor}+1 \leq \frac{\ell+1}{T-\ell}\frac{1}{\ell+1-\lfloor\ell/2\rfloor}$ for even ℓ and $\omega_{\ell+1-\lfloor\ell/2\rfloor}+1 \leq \frac{\ell+1}{T-\ell}\frac{2}{\ell+1-\lfloor\ell/2\rfloor}$ for odd ℓ ." Denote $c_t^u \coloneqq -1 + \frac{\ell+1}{T-\ell}\frac{t'}{\lfloor(\ell+1)/2\rfloor+1}$. If we can show $\omega_t \leq -1 + \frac{\ell+1}{T-\ell}\frac{t'}{\lfloor(\ell+1)/2\rfloor+1} \coloneqq c_t^u$ for $t=t'+\lfloor\ell/2\rfloor$, then it implies " $\omega_{\ell+1-\lfloor\ell/2\rfloor}+1 \leq \frac{\ell+1}{T-\ell}\frac{1}{\ell+1-\lfloor\ell/2\rfloor}$ for even ℓ and $\omega_{\ell+1-\lfloor\ell/2\rfloor}+1 \leq \frac{\ell+1}{T-\ell}\frac{2}{\ell+1-\lfloor\ell/2\rfloor}$ for odd ℓ ." Note that the diagonal entries in $A^{(\ell)}$ are positive while the off-diagonal entries in $A^{(\ell)}$ are negative, then $A_{t',:}^{(\ell)}\vec{\omega}_{(\lfloor\ell/2\rfloor+1):L}$ is increasing in ω_t and decreasing in ω_s for $t=t'+\lfloor\ell/2\rfloor$ and $s\neq t'+\lfloor\ell/2\rfloor$. We only need to show $(A^{(\ell)}c^u)_{\ell-\lfloor\ell/2\rfloor}\geq b_{\ell-\lfloor\ell/2\rfloor}^{(\ell)}$, where $c^u=\begin{bmatrix}c_{\lfloor\ell/2\rfloor+1}^u \cdots c_{\ell-\lfloor\ell/2\rfloor}^u\end{bmatrix}^{\top}$. If ℓ is even, and when $T>\frac{\ell^2+11\ell+2}{8}$,

$$-(A^{(\ell)}c^u)_{\ell-\lfloor\ell/2\rfloor}+b^{(\ell)}_{\ell-\lfloor\ell/2\rfloor}=\frac{\ell(\ell-1)}{T-\ell}-\frac{\ell(\ell+1)}{T-\ell}\left(\frac{\ell}{\ell+2}-\frac{4}{T-\ell}\right)=-\frac{\ell}{T-\ell}\left(\frac{2}{\ell+2}-\frac{\ell+1}{4(T-\ell)}\right)<0.$$

If ℓ is odd, and when $T > \frac{\ell^3 + 13\ell^2 + 7\ell + 3}{8\ell}$ (note that $\frac{\ell^3 + 13\ell^2 + 7\ell + 3}{8\ell} > \frac{\ell^2 + 11\ell + 2}{8}$),

$$-(A^{(\ell)}c^u)_{\ell-\lfloor \ell/2\rfloor} + b^{(\ell)}_{\ell-\lfloor \ell/2\rfloor} = \frac{\ell(\ell-1)}{T-\ell} - \frac{\ell(\ell+1)}{T-\ell} \frac{\ell+1}{\ell+3} + \frac{(\ell+1)^2}{4(T-\ell)^2} = -\frac{1}{T-\ell} \left(\frac{2\ell}{\ell+3} - \frac{(L+2)^2}{4(T-\ell)} \right) < 0.$$

3. ω_t^* for $t > T - \ell$: this is a symmetric case of ω_t^* for $t < \ell$. The proof of ω_t^* for $t > T - \ell$ carries over to this case.

We have verified that the $\vec{\omega}^*$ defined in Eq. (17) satisfies the KKT conditions and the Hessian of the objective function in the quadratic program (25) is positive semi-definite, then $\vec{\omega}^*$ is an optimal solution for the quadratic program (25).

Proof of Theorem 5 In this proof, we first provide an equivalent and simplified form of the objective function $\operatorname{tr}(\mathcal{Z}^{\top}\Sigma_e^{-1}(\Sigma_e - \Gamma(\Gamma^{\top}\Sigma_e^{-1}\Gamma)^{-1}\Gamma^{\top})\Sigma_e^{-1}\mathcal{Z})$ in the integer program (15). We will show that the objective function is a sum of $\ell+1$ convex optimization problems, with each one to be similar as objective function (31) in the proof of Theorem 3. Then we show the optimal solution by using the results in Theorem 3 and Lemma 6.

Conceptually, we show the equivalent form of the objective function in the integer program (15) using the results in Theorem 3 and Lemma 4. We decompose this task into a few steps:

- 1. Calculate $\Theta_1 = \mathcal{Z}^{\top} \Sigma_e^{-1} \mathcal{Z} \in \mathbb{R}^{(\ell+1) \times (\ell+1)}$
- 2. Calculate $\Theta_2 = \mathcal{Z}^{\top} \Sigma_e^{-1} \Gamma (\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1} \mathcal{Z} \in \mathbb{R}^{(\ell+1) \times (\ell+1)}$
 - (a) Calculate $\mathcal{Z}^{\top} \Sigma_e^{-1} \Gamma \in \mathbb{R}^{(\ell+1) \times (N + (T-\ell-1)p)}$
 - (b) Calculate $(\Gamma^{\top} \Sigma_{e}^{-1} \Gamma)^{-1} \in \mathbb{R}^{(N+(t-1-\ell)p) \times (N+(t-1-\ell)p)}$
 - (c) Calculate $(\vec{z}^{(j)})^{\top} \Sigma_e^{-1} \Gamma (\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1} \vec{z}^{(m)} \in R$
 - (d) Calculate $(\vec{z}^{(j)})^{\top} (\Sigma_e^{-1} \Sigma_e^{-1} \Gamma (\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1}) \vec{z}^{(m)} \in R$

where $\vec{z}^{(j)} = \begin{bmatrix} \vec{z}_j^\top & \vec{z}_{j+1}^\top & \cdots & \vec{z}_{T-\ell-1+j}^\top \end{bmatrix}^\top$ for $j = 1, \dots, \ell+1$ and Γ is equal to

$$\Gamma = \begin{bmatrix} \Gamma_{:(\ell+1)} \\ \Gamma_{:(\ell+2)} \\ \vdots \\ \Gamma_{:T} \end{bmatrix} = \begin{bmatrix} \tilde{I} & \tilde{I} & X \\ \tilde{I} & \tilde{I} & X \\ \vdots & & \ddots \\ \tilde{I} & & & \tilde{I} & X \end{bmatrix} = \begin{bmatrix} \tilde{I} & \tilde{X} \\ \tilde{I} & \tilde{X} \\ \vdots & & \ddots \\ \tilde{I} & & & \tilde{X} \end{bmatrix} \in \mathbb{R}^{(N(T-\ell)) \times (N+(T-\ell-1)p)}, \tag{39}$$

where $\tilde{I} = \begin{bmatrix} I_{N-p} \ \mathbf{0}_{N-p,p} \end{bmatrix}^{\top} \in \mathbb{R}^{N \times (N-p)}$, $\mathbf{0}_{N-p,p}$ is a matrix of $0, \ \vec{1} \in \{1\}^N$. Similarly as Theorem 3, $\Sigma_e^{-1} = \operatorname{diag}(\Psi, \Psi, \dots, \Psi) \in \mathbb{R}^{(N(T-\ell)) \times (N(T-\ell))}$, where

$$\Psi = \Sigma_{e_t}^{-1} = (I_N + UU^\top)^{-1} = I_N - U(I_k + U^\top U)^{-1}U^\top \in \mathbb{R}^{N \times N}.$$

For the first step, $\Theta_1 = \mathcal{Z}^{\top} \Sigma_e^{-1} \mathcal{Z}$ has

$$\Theta_{1,jm} = (\vec{z}^{(j)})^{\top} \Sigma_e^{-1} \vec{z}^{(m)} = \sum_{t=1}^{T-\ell} \vec{z}_{j-1+t}^{\top} \Psi \vec{z}_{m-1+t}.$$

For the second step, $\mathcal{Z}^{\top}\Sigma_e^{-1}\Gamma$ has

$$(\vec{z}^{(j)})^{\intercal} \Sigma_e^{-1} \Gamma = (\vec{z}^{(j)})^{\intercal} \begin{bmatrix} \Psi \tilde{I} & \Psi \tilde{X} \\ \Psi \tilde{I} & \Psi \tilde{X} \\ \vdots & & \ddots \\ \Psi \tilde{I} & & \Psi \tilde{X} \end{bmatrix} = \left[\sum_{t=1}^{T-\ell} \vec{z}_{j-1+t}^{\intercal} \Psi \tilde{I} & \vec{z}_j^{\intercal} \Psi \tilde{X} & \cdots & \vec{z}_{T-\ell+j-1}^{\intercal} \Psi \tilde{X} \right] = \left[(\phi^{(j)})^{\intercal} & (\iota^{(j)})^{\intercal} \right],$$

where $(\phi^{(j)})^{\top} = \sum_{t=1}^{T-\ell} \vec{z}_{j-1+t}^{\top} \Psi \tilde{I} = (\sum_{t=1}^{T-\ell} \vec{z}_{j-1+t})^{\top} \Psi \tilde{I}$ and $(\iota^{(j)})^{\top} = \left[\vec{z}_{j}^{\top} \Psi \tilde{X} \cdots \vec{z}_{T-\ell+j-1}^{\top} \Psi \tilde{X}\right]$. Let $\zeta^{(j)} = \frac{1}{T-\ell} \sum_{t=j}^{T-\ell+j-1} \vec{z}_{t} \in \mathbb{R}^{N}$, then $\phi = T(\zeta^{(j)})^{\top} \Psi \tilde{I}$. Note that U and \tilde{X} are orthogonal, we have $\Psi \tilde{X} = \tilde{X}$ and $\tilde{X}^{\top} \Psi \tilde{X} = N \cdot I_{p}$. Then $(\iota^{(j)})^{\top} = \left[\vec{z}_{j}^{\top} \tilde{X} \cdots \vec{z}_{T-\ell+j-1}^{\top} \tilde{X}\right] = \left[N \tilde{\omega}_{j}^{\top} \cdots N \tilde{\omega}_{T-\ell+j-1}^{\top}\right] = N(\tilde{\omega}^{(j)})^{\top} \in \mathbb{R}^{(T-\ell)p}$, where $\tilde{\omega}_{t} = \frac{1}{N} \sum_{i=1}^{N} \tilde{X}_{i} z_{it} \in \mathbb{R}^{p}$.

Next is to calculate $(\Gamma^{\top}\Sigma_e^{-1}\Gamma)^{-1}$. We have similar decomposition as Theorem 3,

$$(\Gamma^{\top}\Sigma_e\Gamma)^{-1} = \begin{bmatrix} \Xi_{11} & \Xi_{12} \\ \Xi_{21} & \Xi_{22} \end{bmatrix} = \begin{bmatrix} M & -M\tilde{M} \\ -\tilde{M}^{\top}M & \bar{M} + \tilde{M}^{\top}M\tilde{M} \end{bmatrix} \in \mathbb{R}^{(N+(T-\ell-1)p))\times (N+(T-\ell-1)p)},$$

where $\Xi_{11} = M$, $\Xi_{12} = -M\tilde{M}$, $\Xi_{21} = -\tilde{M}^{T}M$, $\Xi_{22} = \bar{M} + \tilde{M}^{T}M\tilde{M}$ with

$$\begin{split} M &= \frac{1}{T-\ell} \left(\tilde{I}^\top \Psi \tilde{I} - \tilde{I}^\top \Psi \tilde{X} (\tilde{X}^\top \Psi \tilde{X})^{-1} \tilde{X}^\top \Psi \tilde{I} \right)^{-1} = \frac{1}{T-\ell} \left(\tilde{I}^\top \Psi \tilde{I} - \frac{1}{N} \tilde{X} \tilde{X}^\top \right)^{-1} \in \mathbb{R}^{(N-p) \times (N-p)} \\ \tilde{M} &= \left[\tilde{I}^\top \Psi \tilde{X} (\tilde{X}^\top \Psi \tilde{X})^{-1} \cdots \tilde{I}^\top \Psi \tilde{X} (\tilde{X}^\top \Psi \tilde{X})^{-1} \right] = \left[\frac{1}{N} \tilde{X} \cdots \frac{1}{N} \tilde{X} \right] \in \mathbb{R}^{(N-p) \times ((T-\ell)p)} \\ \bar{M} &= \operatorname{diag}((\tilde{X}^\top \Psi \tilde{X})^{-1}, (\tilde{X}^\top \Psi \tilde{X})^{-1}, \cdots, (\tilde{X}^\top \Psi \tilde{X})^{-1}) = \frac{1}{N} I \in \mathbb{R}^{((T-\ell)p) \times ((T-\ell)p)} \end{split}$$

Similar as Theorem 3, we can simplify M.

$$M = \frac{1}{T - \ell} \left[\left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} + \left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} \tilde{X} \left(N - \tilde{X}^{\top} \left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} \tilde{X} \right)^{-1} \tilde{X}^{\top} \left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} \right]$$

$$\left(\tilde{I}^{\top} \Psi \tilde{I} \right)^{-1} = \left(I_{N-p} - U_{(1)} (I_k + U^{\top} U)^{-1} U_{(1)}^{\top} \right)^{-1} = I_{N-p} + U_{(1)} (I_k + U_{(2)}^{\top} U_{(2)})^{-1} U_{(1)}^{\top},$$

where $U_{(1)} = \begin{bmatrix} u_1 \ u_2 \ \cdots \ u_{N-p} \end{bmatrix}^{\top} \in \mathbb{R}^{(N-p) \times k}$ and $U_{(2)} = \begin{bmatrix} u_{N-p+1} \ \cdots \ u_N \end{bmatrix}^{\top} \in \mathbb{R}^{p \times k}$ $(U = \begin{bmatrix} U_{(1)}^{\top} \ U_{(2)}^{\top} \end{bmatrix}^{\top})$.

Next is to calculate $(\tilde{z}^{(j)})^{\top} \Sigma_e^{-1} \Gamma(\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top} \Sigma_e^{-1} \tilde{z}^{(m)}$ for $1 \leq j, m \leq \ell + 1$. It is equivalent to calculating $[(\phi^{(j)})^{\top} \ (\iota^{(j)})^{\top}] \begin{bmatrix} \Xi_{11} \ \Xi_{12} \ \Xi_{21} \ \Xi_{22} \end{bmatrix} \begin{bmatrix} \phi^{(m)} \ \iota^{(m)} \end{bmatrix} = (\phi^{(j)})^{\top} \Xi_{11} \phi^{(m)} + (\phi^{(j)})^{\top} \Xi_{12} \iota^{(m)} + (\iota^{(j)})^{\top} \Xi_{21} \phi^{(m)} + (\iota^{(j)})^{\top} \Xi_{22} \iota^{(m)},$ where (recall $\zeta^{(j)} = \frac{1}{T-\ell} \sum_{t=j}^{T-\ell+j-1} \tilde{z}_t$ and $\tilde{\omega}^{(j)} = [\tilde{\omega}_j^{\top} \cdots \tilde{\omega}_{T-\ell+j-1}^{\top}]$)

$$\begin{split} &(\phi^{(j)})^{\top}\Xi_{11}\phi^{(m)} = (T-\ell)^{2}(\zeta^{(j)})^{\top}\Psi\tilde{I}M\tilde{I}^{\top}\Psi\zeta^{(m)} \\ &(\phi^{(j)})^{\top}\Xi_{12}\iota^{(m)} = -N(T-\ell)(\zeta^{(j)})^{\top}\Psi\tilde{I}M\tilde{M}\tilde{\omega}^{(m)} = -(T-\ell)(\zeta^{(j)})\Psi\tilde{I}M\tilde{X}\left(\sum_{t=m}^{T-\ell+m-1}\tilde{\omega}_{m-1+t}\right) \\ &(\iota^{(j)})^{\top}\Xi_{21}\phi^{(m)} = -N(T-\ell)(\tilde{\omega}^{(j)})^{\top}\tilde{M}^{\top}M\Psi\tilde{I}\zeta^{(m)} = -(T-\ell)\left(\sum_{t=j}^{T-\ell+j-1}\tilde{\omega}_{t}\right)\tilde{X}^{\top}M\Psi\tilde{I}^{\top}\zeta^{(m)} \\ &(\iota^{(j)})^{\top}\Xi_{22}\iota^{(m)} = N\sum_{t=j}^{T-\ell+j-1}\tilde{\omega}_{t}^{\top}\tilde{\omega}_{t+m-j} + \left(\sum_{t=j}^{T-\ell+j-1}\tilde{\omega}_{t}^{\top}\right)\tilde{X}^{\top}M\tilde{X}\left(\sum_{t=m}^{T-\ell+m-1}\tilde{\omega}_{t}\right) \end{split}$$

 $(\phi^{(j)})^{\top}\Xi_{11}\phi^{(m)}, (\phi^{(j)})^{\top}\Xi_{12}\iota^{(m)}, (\iota^{(j)})^{\top}\Xi_{21}\phi^{(m)}$ and $(\iota^{(j)})^{\top}\Xi_{22}\iota^{(m)}$ can be simplified similar as the proof of Theorem 3 (the definition of notations can be found in the proof of Theorem 3 as well).

$$\begin{split} &(\phi^{(j)})^{\top}\Xi_{11}\phi^{(m)} = (T-\ell)(\zeta^{(j)})^{\top} \left(\Omega + \gamma(NI_{p} - \boldsymbol{\delta})^{-1}\gamma^{\top}\right)\zeta^{(m)} \\ &(\phi^{(j)})^{\top}\Xi_{12}\iota^{(m)} = -(T-\ell)(\zeta^{(j)})^{\top} \left(\gamma(NI_{p} - \boldsymbol{\delta})^{-1}\tilde{X}^{\top}\right)\zeta^{(m)} \\ &(\iota^{(j)})^{\top}\Xi_{21}\phi^{(m)} = -(T-\ell)(\zeta^{(j)})^{\top} \left(\tilde{X}(NI_{p} - \boldsymbol{\delta})^{-1}\gamma^{\top}\right)\zeta^{(m)} \\ &(\iota^{(j)})^{\top}\Xi_{22}\iota^{(m)} = N \sum_{t=j}^{T-\ell+j-1} \tilde{\omega}_{t}^{\top}\tilde{\omega}_{t+m-j} + (T-\ell)(\zeta^{(j)})^{\top} \left(\tilde{X}(NI_{p} - \boldsymbol{\delta})^{-1}\tilde{X}^{\top}\right)\zeta^{(m)} - \frac{N}{T-\ell} \left(\sum_{t=j}^{T-\ell+j-1} \tilde{\omega}_{t}^{\top}\right) \left(\sum_{t=m}^{T-\ell+m-1} \tilde{\omega}_{t}\right) \end{split}$$

Similar as the proof of Theorem 3, we can show

$$\left[(\phi^{(j)})^{\top} (\iota^{(j)})^{\top} \right] \left[\Xi_{11} \Xi_{12} \Xi_{21} \Xi_{22} \right] \left[\phi^{(m)} \right] = N \sum_{t=j}^{T-\ell+j-1} \tilde{\omega}_{t}^{\top} \tilde{\omega}_{t+m-j} - \frac{N}{T-\ell} \left(\sum_{t=j}^{T-\ell+j-1} \tilde{\omega}_{t}^{\top} \right) \left(\sum_{t=m}^{T-\ell+m-1} \tilde{\omega}_{t} \right) + (T-\ell) (\zeta^{(j)})^{\top} \left(I_{N} - U(I_{k} + U^{\top}U)^{-1}U^{\top} \right) \zeta^{(m)} .$$

 $\textbf{Next is to calculate} \ (\vec{z}^{(j)})^\top \, (\Sigma_e^{-1} - \Sigma_e^{-1} \Gamma (\Gamma^\top \Sigma_e^{-1} \Gamma)^{-1} \Gamma^\top \Sigma_e^{-1}) \, \vec{z}^{(m)} \ \textbf{for} \ 1 \leq j, m \leq \ell + 1.$

$$\begin{split} &(\vec{z}^{(j)})^{\top} \boldsymbol{\Sigma}_{e}^{-1} \vec{z}^{(m)} - (\vec{z}^{(j)})^{\top} \boldsymbol{\Sigma}_{e}^{-1} \boldsymbol{\Gamma} (\boldsymbol{\Gamma}^{\top} \boldsymbol{\Sigma}_{e}^{-1} \boldsymbol{\Gamma})^{-1} \boldsymbol{\Gamma}^{\top} \boldsymbol{\Sigma}_{e}^{-1} \vec{z}^{(m)} = \boldsymbol{\Theta}_{1,jm} - \boldsymbol{\Theta}_{2,jm} \\ &= N(T-\ell) - \sum_{t=j}^{T-\ell+j-1} \vec{z}_{t}^{\top} \boldsymbol{U} (I_{k} + \boldsymbol{U}^{\top} \boldsymbol{U})^{-1} \boldsymbol{U}^{\top} \vec{z}_{t+m-j} \\ &- N \sum_{t=j}^{T-\ell+j-1} \tilde{\omega}_{t}^{\top} \tilde{\omega}_{t+m-j} + \frac{N}{T-\ell} \left(\sum_{t=j}^{T-\ell+j-1} \tilde{\omega}_{t}^{\top} \right) \left(\sum_{t=m}^{T-\ell+m-1} \tilde{\omega}_{t} \right) - (T-\ell) (\boldsymbol{\zeta}^{(j)})^{\top} \left(\boldsymbol{I}_{N} - \boldsymbol{U} (\boldsymbol{I}_{k} + \boldsymbol{U}^{\top} \boldsymbol{U})^{-1} \boldsymbol{U}^{\top} \right) \boldsymbol{\zeta}^{(m)} \\ &= N(T-\ell) - \left[(T-\ell) (\boldsymbol{\zeta}^{(j)})^{\top} \boldsymbol{\zeta}^{(m)} - \frac{N}{T-\ell} \left(\sum_{t=j}^{T-\ell+j-1} \boldsymbol{\omega}_{t}^{(1)} \right) \left(\sum_{t=m}^{T-\ell+m-1} \boldsymbol{\omega}_{t}^{(1)} \right) + N \sum_{t=j}^{T-\ell+j-1} \boldsymbol{\omega}_{t}^{(1)} \boldsymbol{\omega}_{t+m-j}^{(1)} \right] \\ &+ N \cdot \sum_{q=2}^{p} \left[- \sum_{t=j}^{T-\ell+j-1} \boldsymbol{\omega}_{t}^{(q)} \boldsymbol{\omega}_{t+m-j}^{(q)} + \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell+j-1} \boldsymbol{\omega}_{t}^{(q)} \right) \left(\sum_{t=m}^{T-\ell+m-1} \boldsymbol{\omega}_{t}^{(q)} \right) \right] \\ &= b^{(j,m)} \\ &+ (T-\ell) (\boldsymbol{\zeta}^{(j)})^{\top} \boldsymbol{U} (\boldsymbol{I}_{k} + \boldsymbol{U}^{\top} \boldsymbol{U})^{-1} \boldsymbol{U}^{\top} \boldsymbol{\zeta}^{(m)} - \sum_{t=j}^{T-\ell+j-1} \vec{z}_{t}^{\top} \boldsymbol{U} (\boldsymbol{I}_{k} + \boldsymbol{U}^{\top} \boldsymbol{U})^{-1} \boldsymbol{U}^{\top} \vec{z}_{t+m-j}, \end{split}$$

where $\omega_t^{(1)} = \frac{1}{N} \sum_{i=1}^N z_{it}$ and $\omega_t^{(q)} = \frac{1}{N} \sum_{i=1}^N X_{i,q-1} z_{it}$ for $q = 2, \dots, p$.

Note that maximizing $\operatorname{tr}(\mathcal{Z}^{\top}\Sigma_e^{-1}(\Sigma_e - \Gamma(\Gamma^{\top}\Sigma_e^{-1}\Gamma)^{-1}\Gamma^{\top})\Sigma_e^{-1}\mathcal{Z})$ is equivalent to minimizing $-\operatorname{tr}(\mathcal{Z}^{\top}\Sigma_e^{-1}(\Sigma_e - \Gamma(\Gamma^{\top}\Sigma_e^{-1}\Gamma)^{-1}\Gamma^{\top})\Sigma_e^{-1}\mathcal{Z})$ and $-\operatorname{tr}(\mathcal{Z}^{\top}\Sigma_e^{-1}(\Sigma_e - \Gamma(\Gamma^{\top}\Sigma_e^{-1}\Gamma)^{-1}\Gamma^{\top})\Sigma_e^{-1}\mathcal{Z}) = -\sum_{j=1}^{\ell+1}(\vec{z}^{(j)})^{\top}(\Sigma_e^{-1} - \Sigma_e^{-1}\Gamma(\Gamma^{\top}\Sigma_e^{-1}\Gamma)^{-1}\Gamma^{\top}\Sigma_e^{-1})\vec{z}^{(j)} = \sum_{j=1}^{\ell+1}(a^{(j,j)} + b^{(j,j)} + c^{(j,j)})$. From Lemma 4, we know

$$\sum_{j=1}^{\ell+1} a^{(j,j)} = N \cdot \sum_{j=1}^{\ell+1} \left[\sum_{t=j}^{T-\ell-1+j} (\omega_t^{(1)})^2 - \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell-1+j} \omega_t^{(1)} \right)^2 + \sum_{t=1}^{T-\ell-1+j} \frac{2(T-\ell-1+2j-2t)}{T-\ell} \omega_t^{(1)} \right].$$

From Lemma 6, $\vec{\omega}^*$ defined in Eq. (17) minimizes $\sum_{j=1}^{\ell+1} a^{(j,j)}$.

$$\sum_{j=1}^{\ell+1} b^{(j,j)} = N \cdot \sum_{j=1}^{\ell+1} \sum_{q=2}^{p} \left[\sum_{t=j}^{T-\ell+j-1} (\omega_t^{(q)})^2 - \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell+j-1} \omega_t^{(q)} \right)^2 \right].$$

From Theorem 3, if $\frac{1}{N} \sum_{i=1}^{N} X_i z_{it} = \mu_X$ for some $\mu_X \in \mathbb{R}^r$ and for all t, then $\sum_{j=1}^{\ell+1} b^{(j,j)}$ is minimized.

$$\begin{split} \sum_{j=1}^{\ell+1} c^{(j,j)} &= \sum_{j=1}^{\ell+1} \left[\sum_{t=j}^{T-\ell+j-1} \vec{z}_t^\top U(I_k + U^\top U)^{-1} U^\top \vec{z}_t - (T-\ell) (\zeta^{(j)})^\top U(I_k + U^\top U)^{-1} U^\top \zeta^{(j)} \right] \\ &= \sum_{j=1}^{\ell+1} \left[\sum_{t=j}^{T-\ell+j-1} \vec{z}_t^\top U(I_k + U^\top U)^{-1} U^\top \vec{z}_t - \frac{1}{T-\ell} \left(\sum_{t=j}^{T-\ell+j-1} \vec{z}_t \right)^\top U(I_k + U^\top U)^{-1} U^\top \left(\sum_{t=j}^{T-\ell+j-1} \vec{z}_t \right) \right] \end{split}$$

From Theorem 3, if $\frac{1}{N} \sum_{i=1}^{N} u_i z_{it} = \mu_U$ for some $\mu_U \in \mathbb{R}^k$ and for all t, then $\sum_{j=1}^{\ell+1} c^{(j,j)}$ is minimized. Hence, if the optimal design satisfies

$$\frac{1}{N} \sum_{i=1}^{N} z_{it} = \omega_t^*, \qquad \frac{1}{N} \sum_{i=1}^{N} X_i z_{it} = \mu_X, \quad \frac{1}{N} \sum_{i=1}^{N} u_i z_{it} = \mu_U, \quad \text{for all } t$$
(40)

for some $\mu_X \in \mathbb{R}^r$ and $\mu_U \in \mathbb{R}^k$. Then $\sum_{j=1}^{\ell+1} a^{(j,j)}$, $\sum_{j=1}^{\ell+1} b^{(j,j)}$ and $\sum_{j=1}^{\ell+1} c^{(j,j)}$ are all minimized and therefore $\sum_{j=1}^{\ell+1} (a^{(j,j)} + b^{(j,j)} + c^{(j,j)})$ is minimized, which is equivalent to $\operatorname{tr}(\mathcal{Z}^{\top} \Sigma_e^{-1} (\Sigma_e - \Gamma(\Gamma^{\top} \Sigma_e^{-1} \Gamma)^{-1} \Gamma^{\top}) \Sigma_e^{-1} \mathcal{Z})$ being maximized. \square

Appendix F: Additional Empirical Results

F.1. Additional Empirical Datasets

Offering extra reward points or discounts (if they spend a certain amount of money on a shopping trip). Increasing the marginal redemption rates if customers redeem more reward points. Offering personalized promotions on the categories that customers shop less frequently.

Our analysis focuses on frequent shoppers because they are more familiar and tend to pay more attention to the reward or loyalty program's changes. We define frequent shoppers as the households with expenditure in at least half of the weeks in the data set, that is, more than 48 of the 97 total weeks. There are 7,130 frequent households. We want to study if a new approach can increase frequent households expenditure.

Home Medical Visits. This data set has 40,079 records of home medical visits from Jan 2016 to Dec 2018 in Barcelona Spain, which is publicly available on Kaggle.¹¹ We aggregate records by city and week, calculate visit rates by taking the moving average of numbers of visits, and get a panel of 102 cities over 144 weeks. We consider home medical visit rates instead of the raw number of visits for two reasons:

- 1. Home visits can be naturally modeled as a Poisson distribution. We can calculate the visit rate by averaging the number of home visits in the past m weeks, which is an unbiased and consistent estimator for the rate (we use m = 16 here). It is natural and reasonable to measure the impact of environmental policies by the impact on the visit rate.
- 2. We are interested in the environmental policies that can improve public health, i.e., reduce the number of home visits. If an entry has value zero, the value could not be further reduced. However, if we study visit rates, most entries are positive, which can potentially be reduced.

Transactions from a Large Grocery Store. This data set contains 17,880,248 transactions between May 2005 and May 2007. We aggregate transactions by household and week and get a panel of 7,130 frequent households over 97 weeks, where frequent households had expenditure in more than 50% weeks.

¹¹ This data set can be downloaded at https://www.kaggle.com/HackandHealth/home-medical-visits-healthcare

We are interested in studying if some new approaches could incentivize households to increase their expenditure in this grocery store. These new approaches could be: 1. increasing the number of queues (Lu et al. 2013); 2. offering extra reward points or discounts; 3. increasing the marginal redemption rates if customers redeem more reward points; 4. offering personalized promotions on the categories that customers shop less frequently; 5. designing a better search engine for this grocery store's website and app (lEcuyer et al. 2017).

F.2. Additional Experiment Setups

Specification. In this section, the results are calculated based on 100 randomly selected blocks divided into the historical control and synthetic experimental data. For the home medical visit data, we choose the treatment effect at $\tau = -0.05$ in the experiment with direct effect only. Here the median value is 1.3125, so the treatment effect is about 4% of the median visit rate. The lowest visit rate is 0.0625, so the τ we choose allows the treated outcome to stay positive. For the grocery data, we choose the treatment effect at $\tau = 10$ in the experiment with direct effect only. Note that the median expenditure is 43.54 so the scale of the treatment effect is about 23% of the median expenditure. In the experiment with carryover effects, we consider the treatment can affect the current period and four periods in the future (equivalent to $\ell = 4$) and choose the treatment effects at $[\tau_1, \tau_2, \tau_3, \tau_4, \tau_5] = [-0.03, -0.01, -0.005, -0.002, -0.001]$ (so the treated outcome stays positive) for the home medical visit data and $[\tau_1, \tau_2, \tau_3, \tau_4, \tau_5] = [6, 4, 2, 1, 0.5]$ for the grocery data.

Choice of Estimation Method. In practice, we need to select an estimation method to estimate the treatment effect on the experimental data. We can leverage historical control data and use the procedure in Remark 14 to find the best estimation method. Then we use this method to estimate the treatment effect on the synthetic experimental data. The results are presented in the row "Hist Winner" in Tables 6 and 8.

REMARK 17. One possible reason for the BLUE estimator to perform worse than LRME on the home visit data is that the visit rates are strongly time-series correlated by construction. In GLS, we calculate errors' covariance matrix based on the assumption that errors are time-series uncorrelated and use the inverse of this covariance matrix as the weighting matrix. This assumption allows us to estimate fewer parameters $(N \times N)$ rather than $NT \times NT$ in the covariance matrix. Hence feasible GLS is less efficient when the assumption on time-series uncorrelation is violated.

 $^{^{12}}$ We also estimate the full errors' covariance matrix with $NT \times NT$ parameters that do not rely on the time-series uncorrelated assumption. However, the corresponding RMSE is much larger than the RMSE presented in Table 6.

| | (N,T) | (25,10) | (25,20) | (50,10) | (50,20) | | (N,T) | (25,10) | (25,20) | (50,10) | (50, |
|----------------|----------------------|---------|---------|---------|---------|----------------|---------------------|---------|---------|---------|------|
| OLS | $ Z_{ m BA} $ | 0.1500 | 0.2278 | 0.1351 | 0.2266 | OLS | Z_{OPT} | 0.1389 | 0.2535 | 0.0681 | 0.15 |
| | $Z_{ m FF}$ | 0.4090 | 0.7363 | 0.3855 | 0.5681 | | $Z_{K=2}$ | 0.0906 | 0.1116 | 0.0568 | 0.09 |
| | $Z_{\mathrm{FF+BA}}$ | 0.2655 | 0.4708 | 0.1436 | 0.3169 | | $Z_{K=3}$ | 0.1011 | 0.1126 | 0.0557 | 0.07 |
| | $Z_{ m OPT}$ | 0.1389 | 0.2535 | 0.0681 | 0.1577 | | $Z_{\mathrm{OPT+}}$ | 0.0552 | 0.0817 | 0.0380 | 0.06 |
| GLS | $ Z_{\mathrm{BA}} $ | 0.0653 | 0.0870 | 0.0695 | 0.0873 | GLS | Z_{OPT} | 0.0582 | 0.1031 | 0.0321 | 0.06 |
| | $Z_{ m FF}$ | 0.1846 | 0.1677 | 0.1653 | 0.1765 | | $Z_{K=2}$ | 0.0423 | 0.0536 | 0.0285 | 0.04 |
| | $Z_{\rm FF+BA}$ | 0.2033 | 0.3531 | 0.1072 | 0.2410 | | $Z_{K=3}$ | 0.0474 | 0.0486 | 0.0285 | 0.03 |
| | $Z_{ m OPT}$ | 0.0582 | 0.1031 | 0.0321 | 0.0645 | | $Z_{\mathrm{OPT+}}$ | 0.0284 | 0.0381 | 0.0190 | 0.03 |
| LRME | $ Z_{ m BA} $ | 0.0724 | 0.0863 | 0.0634 | 0.0703 | LRME | Z_{OPT} | 0.0292 | 0.0306 | 0.0218 | 0.02 |
| | $Z_{ m FF}$ | 0.2116 | 0.3308 | 0.1027 | 0.0986 | | $Z_{K=2}$ | 0.0268 | 0.0338 | 0.0220 | 0.02 |
| | $Z_{\mathrm{FF+BA}}$ | 0.1877 | 0.3924 | 0.0744 | 0.2101 | | $Z_{K=3}$ | 0.0265 | 0.0287 | 0.0224 | 0.02 |
| | $Z_{ m OPT}$ | 0.0292 | 0.0306 | 0.0218 | 0.0223 | | $Z_{\mathrm{OPT+}}$ | 0.0261 | 0.0279 | 0.0214 | 0.02 |
| Hist Winner | $Z_{ m OPT}$ | 0.0323 | 0.0308 | 0.0240 | 0.0237 | Hist Winner | $Z_{\mathrm{OPT+}}$ | 0.0264 | 0.0299 | 0.0219 | 0.02 |

Table 6 Home medical visit data: This table compares the RMSE based on m = 100 randomly sampled blocks for benchmark treatment designs, $Z_{\rm OPT}$, $Z_{K=2}$, $Z_{K=3}$, $Z_{\rm OPT+}$ and "Hist Winner" (the estimation method with minimax estimation error on historical matrices). In the synthetic experiments, the intervention only has *direct* treatment effect and we choose the treatment effect at $\tau = -0.05$. The left table shows the linear staggered design $Z_{\rm OPT}$ outperforms benchmark treatment designs. The right table shows we can find a better treatment design $Z_{\rm OPT+}$ via historical data and our data-driven local search algorithm. Both tables show LRME is the best estimation method on the home medical visit data and the optimal estimation method found on historical data has similar performance as LRME.

| | (N,T) | (25,10) | (25,20) | (50,10) | (50,20) |
|------|---|--|-------------------------|-------------------------|-------------------------|
| OLS | $\left egin{array}{c} Z_{ m OPT} \ Z_{ m OPT-CO} \end{array} ight $ | $\begin{vmatrix} 0.1019 \\ 0.0938 \end{vmatrix}$ | $0.1476 \\ 0.1360$ | $0.0612 \\ 0.0583$ | 0.0968 0.0898 |
| GLS | $\left egin{array}{c} Z_{ m OPT-CO} \end{array} ight $ | $\begin{vmatrix} 0.0640 \\ 0.0593 \end{vmatrix}$ | $0.0752 \\ 0.0591$ | 0.0432 0.0405 | $0.0482 \\ 0.0423$ |
| LRME | $ Z_{ m OPT} $ | 0.0501 0.0478 | 0.0407 0.0405 | $0.0435 \\ 0.0418$ | 0.0410 0.0375 |

Table 7 Home medical visit data: This table compares the RMSE based on m=100 randomly sampled blocks for $Z_{\rm OPT}$ and $Z_{\rm OPT-CO}$. In the synthetic experiments, the intervention has direct and carryover treatment effects and we choose the treatment effects at $[\tau_1, \tau_2, \tau_3, \tau_4, \tau_5] = [-0.03, -0.01, -0.005, -0.002, -0.001]$. This table shows the nonlinear staggered design $Z_{\rm OPT-CO}$ outperforms the linear staggered design $Z_{\rm OPT}$.

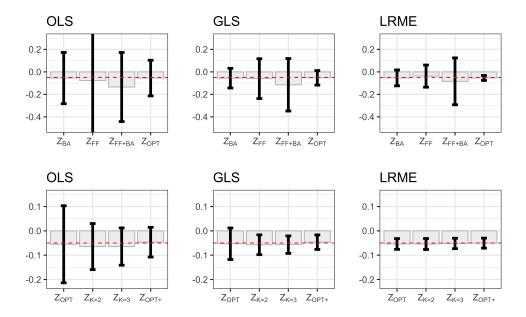


Figure 5 Home medical visit data: This figure compares the average estimated treatment effect $\hat{\tau}_Z$ and its standard deviation $\sqrt{\mathrm{Var}(\hat{\tau}_Z)}$ based on m=100 randomly sampled blocks for benchmark designs, $Z_{\mathrm{OPT}}, Z_{K=2}, Z_{K=3}$ and $Z_{\mathrm{OPT+}}$. In the synthetic experiments, N=50, T=20, the intervention only has direct treatment effect and we choose the treatment effect at $\tau=-0.05$. The height of the bar shows $\hat{\tau}_Z$, while the error bar indicates the standard deviation $\sqrt{\mathrm{Var}(\hat{\tau}_Z)}$. The red dash line indicates the true value of $\tau=-0.05$. Note that figures in the second row have a different y-axis scale due to superior performance of $Z_{\mathrm{OPT}}, Z_{K=2}, Z_{K=3}$ and $Z_{\mathrm{OPT+}}$ over benchmark treatment designs. The bias of various treatment designs is similar while Z_{OPT} has much smaller variance compared with benchmark designs and the treatment designs found using historical data, such as $Z_{\mathrm{OPT+}}$ from Algorithm 1, has smaller variance compared with Z_{OPT} .

| | (N,T) | (50,10) | (50,20) | (100,10) | (100,20) | | _ | (N,T) | (50,10) | (50,20) | (100,10) | (100 |
|----------------|----------------------|---------|---------|----------|----------|------------|----|---------------------|---------|---------|----------|-------|
| OLS | $ Z_{ m BA} $ | 35.1941 | 22.7694 | 23.7481 | 21.7461 | OLS | | $Z_{ m OPT}$ | 4.9232 | 4.1762 | 3.6979 | 2.70 |
| | $Z_{ m FF}$ | 19.4140 | 18.7226 | 19.9693 | 21.7770 | | | $Z_{K=2}$ | 5.4849 | 4.1601 | 3.4513 | 2.50 |
| | $Z_{\mathrm{FF+BA}}$ | 4.6860 | 4.2707 | 3.7873 | 2.9318 | | | $Z_{K=3}$ | 5.1240 | 4.0824 | 3.7250 | 2.54 |
| | $Z_{ m OPT}$ | 4.9232 | 4.1762 | 3.6979 | 2.7044 | | | $Z_{\mathrm{OPT+}}$ | 4.8153 | 2.6892 | 3.3844 | 2.48 |
| GLS | Z_{BA} | 32.0785 | 16.0023 | 19.8024 | 17.5874 | GLS | | $Z_{ m OPT}$ | 2.4448 | 1.8763 | 1.8044 | 1.48 |
| | $Z_{ m FF}$ | 15.7801 | 15.8220 | 18.9532 | 21.9945 | | | $Z_{K=2}$ | 2.6511 | 1.8308 | 1.7126 | 1.44 |
| | $Z_{\rm FF+BA}$ | 2.6314 | 2.3725 | 2.1481 | 1.6479 | | | $Z_{K=3}$ | 2.4610 | 1.8500 | 1.7607 | 1.39 |
| | $Z_{ m OPT}$ | 2.4448 | 1.8763 | 1.8044 | 1.4877 | | | $Z_{\mathrm{OPT+}}$ | 2.0462 | 1.5986 | 1.6463 | 1.30 |
| LRME | Z_{BA} | 50.2295 | 23.5581 | 26.9896 | 22.7613 | LRM | IE | $Z_{ m OPT}$ | 4.9026 | 3.9167 | 3.8415 | 2.63 |
| | $Z_{ m FF}$ | 25.6949 | 21.4246 | 22.1375 | 22.9238 | | | $Z_{K=2}$ | 5.1930 | 3.9344 | 3.4794 | 2.426 |
| | $Z_{\rm FF+BA}$ | 5.5279 | 4.4794 | 4.2237 | 2.8163 | | | $Z_{K=3}$ | 5.0724 | 3.8695 | 3.7010 | 2.543 |
| | $Z_{ m OPT}$ | 4.9026 | 3.9167 | 3.8415 | 2.6328 | | | $Z_{\mathrm{OPT+}}$ | 4.7684 | 3.3779 | 3.5350 | 2.28 |
| Hist Winner | $Z_{ m OPT}$ | 3.1680 | 2.1187 | 2.1495 | 1.7049 | His Win | | $Z_{ m OPT+}$ | 3.1691 | 1.8382 | 1.7999 | 1.37 |

Table 8 Grocery data: This table compares the RMSE based on m = 100 randomly sampled blocks for benchmark treatment designs, $Z_{\rm OPT}$, $Z_{K=2}$, $Z_{K=3}$, $Z_{\rm OPT+}$ and "Hist Winner" (the estimation method with minimax estimation error on historical matrices). In the synthetic experiments, the intervention only has direct treatment effect and we choose the treatment effect at $\tau = 10$. The left table shows the linear staggered design $Z_{\rm OPT}$ outperforms benchmark treatment designs. The right table shows we can find a better treatment design $Z_{\rm OPT+}$ via historical data and our data-driven local search algorithm. Both tables show GLS is the best estimation method and the optimal estimation method found on historical data has similar performance as GLS.

| | (N,T) | (50,10) | (50,20) | (100,10) | (100,20) |
|------|---|--|-------------------------|-------------------------|-------------------------|
| OLS | $egin{array}{c} Z_{ m OPT-CO} \end{array}$ | $\begin{vmatrix} 10.0771 \\ 9.6095 \end{vmatrix}$ | 8.9105 8.9056 | 6.9050 6.8567 | 6.0699 5.8711 |
| GLS | $ig _{Z_{ m OPT-CO}}$ | 7.1545 6.9175 | 6.2967 6.1221 | 5.1317 4.9333 | 4.5215 4.4452 |
| LRME | $egin{array}{c} Z_{ m OPT} \ Z_{ m OPT-CO} \end{array}$ | $\begin{vmatrix} 12.9216 \\ 12.3249 \end{vmatrix}$ | 10.7967 10.7979 | 10.0112 9.3749 | 7.9319 7.7643 |

Table 9 Grocery data: This table compares the RMSE based on m = 100 randomly sampled blocks for Z_{OPT} and $Z_{\text{OPT-CO}}$. In the synthetic experiments, the intervention has *direct* and *carryover* treatment effects and we choose the treatment effects at $[\tau_1, \tau_2, \tau_3, \tau_4, \tau_5] = [6, 4, 2, 1, 0.5]$. This table shows the *nonlinear* staggered design $Z_{\text{OPT-CO}}$ outperforms the linear staggered design Z_{OPT} .

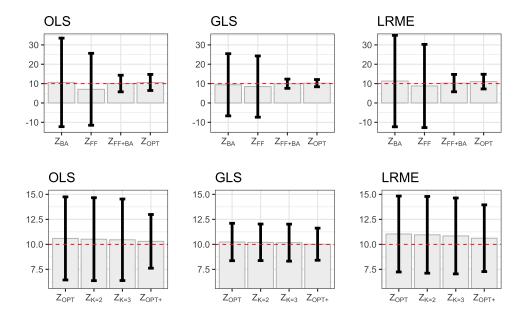


Figure 6 Grocery data: This figure compares the average estimated treatment effect $\hat{\tau}_Z$ and its standard deviation $\sqrt{\operatorname{Var}(\hat{\tau}_Z)}$ based on m=100 randomly sampled blocks for benchmark designs, $Z_{\mathrm{OPT}}, Z_{K=2}, Z_{K=3}$ and $Z_{\mathrm{OPT}+}$. In the synthetic experiments, N=50, T=20, the intervention only has direct treatment effect and we choose the treatment effect at $\tau=10$. The height of the bar shows $\hat{\tau}_Z$, while the error bar indicates the standard deviation $\sqrt{\operatorname{Var}(\hat{\tau}_Z)}$. The red dash line indicates the true value of $\tau=10$. Note that figures in the second row have a different y-axis scale due to superior performance of $Z_{\mathrm{OPT}}, Z_{K=2}, Z_{K=3}$ and $Z_{\mathrm{OPT}+}$ over benchmark treatment designs. The bias of various treatment designs is similar while Z_{OPT} has much smaller variance compared with benchmark designs and the treatment designs found using historical data, such as $Z_{\mathrm{OPT}+}$ from Algorithm 1, has smaller variance compared with Z_{OPT} .