

UC Berkeley

UC Berkeley Previously Published Works

Title

The Augmented Synthetic Control Method

Permalink

<https://escholarship.org/uc/item/27b8831v>

Authors

Ben-Michael, Eli

Feller, Avi

Rothstein, Jesse

Publication Date

2021

DOI

10.2139/ssrn.3861414

Peer reviewed

The Augmented Synthetic Control Method*

Eli Ben-Michael, Avi Feller, and Jesse Rothstein

UC Berkeley

July 2020

Abstract

The synthetic control method (SCM) is a popular approach for estimating the impact of a treatment on a single unit in panel data settings. The “synthetic control” is a weighted average of control units that balances the treated unit’s pre-treatment outcomes as closely as possible. A critical feature of the original proposal is to use SCM only when the fit on pre-treatment outcomes is excellent. We propose Augmented SCM as an extension of SCM to settings where such pre-treatment fit is infeasible. Analogous to bias correction for inexact matching, Augmented SCM uses an outcome model to estimate the bias due to imperfect pre-treatment fit and then de-biases the original SCM estimate. Our main proposal, which uses ridge regression as the outcome model, directly controls pre-treatment fit while minimizing extrapolation from the convex hull. This estimator can also be expressed as a solution to a modified synthetic controls problem that allows negative weights on some donor units. We bound the estimation error of this approach under different data generating processes, including a linear factor model, and show how regularization helps to avoid over-fitting to noise. We demonstrate gains from Augmented SCM with extensive simulation studies and apply this framework to estimate the impact of the 2012 Kansas tax cuts on economic growth. We implement the proposed method in the new `augsynth` R package.

*email: afeller@berkeley.edu. We thank Alberto Abadie, Josh Angrist, Matias Cattaneo, Alex D’Amour, Peng Ding, Erin Hartman, Chad Hazlett, Steve Howard, Guido Imbens, Brian Jacob, Pat Kline, Caleb Miles, Luke Miratrix, Sam Pimentel, Fredrik Sävje, Jas Sekhon, Jake Soloff, Panos Toulis, Stefan Wager, Yiqing Xu, Alan Zaslavsky, and Xiang Zhou for thoughtful comments and discussion, as well as seminar participants at Stanford, UC Berkeley, UNC, the 2018 Atlantic Causal Inference Conference, COMPIE 2018, and the 2018 Polmeth Conference. We also thank editors and referees for constructive feedback.

1 Introduction

The *synthetic control method* (SCM) is a popular approach for estimating the impact of a treatment on a single unit in panel data settings with a modest number of control units and with many pre-treatment periods (Abadie and Gardeazabal, 2003; Abadie et al., 2010, 2015). The idea is to construct a weighted average of control units, known as a synthetic control, that matches the treated unit’s pre-treatment outcomes. The estimated impact is then the difference in post-treatment outcomes between the treated unit and the synthetic control. SCM has been widely applied — the main SCM papers have over 4,000 citations — and has been called “arguably the most important innovation in the policy evaluation literature in the last 15 years” (Athey and Imbens, 2017).

A critical feature of the original proposal, not always followed in practice, is to use SCM only when the synthetic control’s pre-treatment outcomes closely match the pre-treatment outcomes for the treated unit (Abadie et al., 2015). When it is not possible to construct a synthetic control that fits pre-treatment outcomes well, the original papers advise against using SCM. At that point, researchers often fall back to linear regression. This allows better (often perfect) pre-treatment fit, but does so by applying negative weights to some control units, extrapolating outside the support of the data.

We propose the *augmented synthetic control method* (ASCM) as a middle ground in settings where excellent pre-treatment fit using SCM alone is not feasible. Analogous to bias correction for inexact matching (Rubin, 1973; Abadie and Imbens, 2011), ASCM begins with the original SCM estimate, uses an outcome model to estimate the bias due to imperfect pre-treatment fit, and then uses this to de-bias the estimate. If pre-treatment fit is good, the estimated bias will be small, and the SCM and ASCM estimates will be similar. Otherwise, the estimates will diverge, and ASCM will rely more heavily on extrapolation.

Our primary proposal is to augment SCM with a ridge regression model, which we call *Ridge ASCM*. We show that, like SCM, the Ridge ASCM estimator can be written as a weighted average of the control unit outcomes. We also show that Ridge ASCM weights can be written as the solution to a modified synthetic controls problem, targeting the same imbalance metric as traditional SCM. However, where SCM weights are always non-negative, Ridge ASCM admits negative weights, using extrapolation to improve pre-treatment fit. The regularization parameter in Ridge ASCM directly parameterizes the level of extrapolation by penalizing the distance from SCM weights. By contrast, (ridge) regression alone, which can also be written as a modified synthetic controls problem with possibly negative weights, allows for arbitrary extrapolation and possibly unchecked extrapolation bias.

We relate Ridge ASCM’s improved pre-treatment fit to a finite sample bound on estimation error under several data generating processes, including a simple linear model and the linear factor model often invoked in this setting (Abadie et al., 2010). Under a linear model, improving pre-treatment fit directly reduces bias, and the Ridge ASCM penalty term negotiates a bias-variance

trade-off. Under a latent factor model, improving pre-treatment fit again reduces bias, though there is now a risk of over-fitting. The penalty term also directly parameterizes this trade-off. Thus, choosing the hyperparameter will be important for practice; we propose a cross-validation procedure in Section 5.3.

Finally, we also describe how the Augmented SCM approach can be extended to incorporate auxiliary covariates other than pre-treatment outcomes. We first propose to include the auxiliary covariates in parallel to the lagged outcomes in both the SCM and outcome models. We also propose an alternative when there are relatively few covariates, extending a suggestion from Doudchenko and Imbens (2017): first residualize both the pre- and post-treatment outcomes against the auxiliary covariates, then fit Ridge ASCM on the residualized outcome series. We show that this controls the estimation error under a linear factor model with auxiliary covariates.

We demonstrate the properties of Augmented SCM both via calibrated simulation studies and by using it to examine the effect of an aggressive tax cut in Kansas in 2012 on economic output, finding a substantial negative effect. Overall, we see large gains from ASCM relative to alternative estimators, especially under model mis-specification, in terms of both bias and root mean squared error. We implement the proposed methodology in the `augsynth` package for R, available at <https://github.com/ebenmichael/augsynth>.

The paper proceeds as follows. Section 1.1 briefly reviews related work. Section 2 introduces notation and the SCM estimator. Section 3 gives an overview of Augmented SCM. Section 4 gives key numerical results for Ridge ASCM. Section 5 bounds the estimation error under a linear model and a linear factor model, the standard setting for SCM. Section 6 extends the ASCM framework to incorporate auxiliary covariates and alternative outcome models. Section 7 reports on extensive simulation studies as well as the application to the Kansas tax cuts. Finally, Section 8 discusses some possible directions for further research. The appendix includes all of the proofs, as well as additional derivations and technical discussion, including a discussion of inference.

1.1 Related work

SCM was introduced by Abadie and Gardeazabal (2003) and Abadie et al. (2010, 2015) and is the subject of an extensive methodological literature; see Abadie (2019) and Samartsidis et al. (2019) for recent reviews. We briefly highlight some relevant aspects of this literature.

A first set of papers adapts the original SCM proposal to allow for more robust estimation while retaining the simplex constraint on the weights. Robbins et al. (2017); Doudchenko and Imbens (2017); Abadie and L'Hour (2018); Minard and Waddell (2018) incorporate a penalty on the weights into the SCM optimization problem, building on a suggestion in Abadie et al. (2015). Gobillon and Magnac (2016) and Hazlett and Xu (2018) explore dimension reduction strategies and other data transformations that can improve the performance of the subsequent estimator. See also Bilinski and Hatfield (2020), who propose a cross-validation procedure for selecting the number of lagged

outcomes to include.

A second set of papers relaxes various constraints imposed in the original SCM problem. In particular, [Doudchenko and Imbens \(2017\)](#) relax the SCM restriction that control unit weights be non-negative, arguing that there are many settings in which negative weights would be desirable. [Amjad et al. \(2018\)](#) propose an interesting variant that combines negative weights with a pre-processing step. [Powell \(2018\)](#) instead allows for extrapolation via a Frisch-Waugh-Lovell-style projection, which similarly generalizes the typical SCM setting. [Doudchenko and Imbens \(2017\)](#) and [Ferman and Pinto \(2018\)](#) both propose to incorporate an intercept into the SCM problem, which we discuss in Section 3.2. There have also been several other proposals to reduce bias in SCM, developed independently and contemporaneously with ours. [Abadie and L'Hour \(2018\)](#) also propose bias correcting SCM using regression, but focus on bias due to interpolation rather than poor pre-treatment fit. [Kellogg et al. \(2020\)](#) propose using a weighted average of SCM and matching, trading off interpolation and extrapolation bias. Finally, [Arkhangelsky et al. \(2019\)](#) propose the *Synthetic Difference-in-Differences* estimator, which can be seen as a special case of our proposal with a constrained outcome regression.

Finally, there have been several recent proposals to use outcome modeling rather than SCM-style weighting in this setting. These include the matrix completion method in [Athey et al. \(2017\)](#), the generalized synthetic control method in [Xu \(2017\)](#), and the combined approaches in [Hsiao et al. \(2018\)](#). We explore the performance of select methods, both in isolation and in combination with SCM, in Section 7.1.

2 Overview of the Synthetic Control Method

2.1 Notation and setup

We consider the canonical SCM panel data setting with $i = 1, \dots, N$ units observed for $t = 1, \dots, T$ time periods; for the theoretical discussion below, we will consider both N and T to be fixed. Let W_i be an indicator that unit i is treated at time $T_0 < T$ where units with $W_i = 0$ never receive the treatment. We restrict our attention to the case where a single unit receives treatment, and follow the convention that this is the first one, $W_1 = 1$; see [Ben-Michael et al. \(2019\)](#) for an extension to multiple treated units. The remaining $N_0 = N - 1$ units are possible controls, often referred to as *donor units* in the SCM context. To simplify notation, we limit to one post-treatment observation, $T = T_0 + 1$, though our results are easily extended to larger T .

We adopt the potential outcomes framework ([Neyman, 1923](#); [Rubin, 1974](#)) and invoke SUTVA, which assumes a well-defined treatment and excludes interference between units ([Rubin, 1980](#)); the potential outcomes for unit i in period t under control and treatment are $Y_{it}(0)$ and $Y_{it}(1)$,

respectively. The observed outcomes are then:

$$Y_{it} = \begin{cases} Y_{it}(0) & \text{if } W_i = 0 \text{ or } t \leq T_0 \\ Y_{it}(1) & \text{if } W_i = 1 \text{ and } t > T_0. \end{cases} \quad (1)$$

We next assume that control potential outcomes are generated as a fixed component m_{it} plus mean zero, additive noise ε_{it} drawn from some distribution $P(\cdot)$,

$$Y_{it}(0) = m_{it} + \varepsilon_{it}.$$

The treated potential outcome is then $Y_{it}(1) = Y_{it}(0) + \tau_{it}$, where the treatment effects τ_{it} are the key estimands, and are fixed parameters. The treatment effect of interest is thus $\tau = \tau_{1T} = Y_{1T}(1) - Y_{1T}(0)$.

In Section 5, we consider special cases where m_{it} is a linear function of lagged outcomes or where m_{it} is a linear factor model; in the Appendix, we also consider the case where m_{it} is a linear model with Lipschitz deviations from linearity. We make two assumptions about the distribution of the noise terms ε_{it} . First, we assume that treatment assignment W_i is ignorable given m_{it} ; specifically that the noise terms in the post-treatment time periods $\varepsilon_T = (\varepsilon_{1T}, \dots, \varepsilon_{NT})$ are mean-zero and uncorrelated with treatment assignment,

$$\mathbb{E}_{\varepsilon_T} [W_i \varepsilon_{iT}] = \mathbb{E}_{\varepsilon_T} [(1 - W_i) \varepsilon_{iT}] = \mathbb{E}_{\varepsilon_T} [\varepsilon_{iT}] = 0, \quad (2)$$

where the expectation is taken with respect to the noise term at the post-treatment time T , ε_T . As a result, the noise terms for the treated and control units do not systematically deviate from each other. We discuss this in more detail in the context of our application in Section 7. Second, for the theoretical results in Sections 5 and 6, we assume that the error terms ε_{it} are independent (across units and over time) sub-Gaussian random variables with scale parameter σ . See, for example, Chernozhukov et al. (2019) for an extended discussion of this general setup in panel data settings.

To emphasize that pre-treatment outcomes serve as covariates in SCM, we use X_{it} , for $t \leq T_0$, to represent pre-treatment outcomes; we use the terms *pre-treatment fit* and *covariate balance* interchangeably. With some abuse of notation, we use \mathbf{X}_0 to represent the N_0 -by- T_0 matrix of control unit pre-treatment outcomes and \mathbf{Y}_{0T} for the N_0 -vector of control unit outcomes in period T . With only one treated unit, Y_{1T} is a scalar, and \mathbf{X}_1 is a T_0 -row vector of treated unit pre-treatment outcomes. The data structure is then:

$$\begin{pmatrix} Y_{11} & Y_{12} & \dots & Y_{1T_0} & Y_{1T} \\ Y_{21} & Y_{22} & \dots & Y_{2T_0} & Y_{2T} \\ \vdots & & & & \vdots \\ Y_{N1} & Y_{N2} & \dots & Y_{NT_0} & Y_{NT} \end{pmatrix} \equiv \left(\begin{array}{cccc|c} X_{11} & X_{12} & \dots & X_{1T_0} & Y_{1T} \\ X_{21} & X_{22} & \dots & X_{2T_0} & Y_{2T} \\ \vdots & & & & \vdots \\ X_{N1} & X_{N2} & \dots & X_{NT_0} & Y_{NT} \end{array} \right) \equiv \left(\begin{array}{c|c} \mathbf{X}_1 & Y_{1T} \\ \hline \mathbf{X}_0 & \mathbf{Y}_{0T} \end{array} \right) \quad (3)$$

pre-treatment outcomes

2.2 Synthetic Control Method

The Synthetic Control Method imputes the missing potential outcome for the treated unit, $Y_{1T}(0)$, as a weighted average of the control outcomes, $\mathbf{Y}'_{0T}\boldsymbol{\gamma}$ (Abadie and Gardeazabal, 2003; Abadie et al., 2010, 2015). Weights are chosen to balance pre-treatment outcomes and possibly other covariates. We consider a version of SCM that chooses weights $\boldsymbol{\gamma}$ as a solution to the constrained optimization problem:

$$\begin{aligned} \min_{\boldsymbol{\gamma}} \quad & \|\mathbf{V}_{\mathbf{x}}^{1/2}(\mathbf{X}_1 - \mathbf{X}'_0\boldsymbol{\gamma})\|_2^2 + \zeta \sum_{W_i=0} f(\gamma_i) \\ \text{subject to} \quad & \sum_{W_i=0} \gamma_i = 1 \\ & \gamma_i \geq 0 \quad i : W_i = 0 \end{aligned} \quad (4)$$

where the constraints limit $\boldsymbol{\gamma}$ to the simplex $\Delta^{N_0} = \{\boldsymbol{\gamma} \in \mathbb{R}^{N_0} \mid \gamma_i \geq 0 \forall i, \sum_i \gamma_i = 1\}$, and where $\mathbf{V}_{\mathbf{x}} \in \mathbb{R}^{T_0 \times T_0}$ is a symmetric importance matrix and $\|\mathbf{V}_{\mathbf{x}}^{1/2}(\mathbf{X}_1 - \mathbf{X}'_0\boldsymbol{\gamma})\|_2^2 \equiv (\mathbf{X}_1 - \mathbf{X}'_0\boldsymbol{\gamma})' \mathbf{V}_{\mathbf{x}} (\mathbf{X}_1 - \mathbf{X}'_0\boldsymbol{\gamma})$ is the 2-norm on \mathbb{R}^{T_0} after applying $\mathbf{V}_{\mathbf{x}}^{1/2}$ as a linear transformation. To simplify the exposition and notation below, we will generally take $\mathbf{V}_{\mathbf{x}}$ to be the identity matrix. The simplex constraint in Equation (4) ensures that the weights will be sparse and non-negative; Abadie et al. (2010, 2015) argue that enforcing this constraint is important for preserving interpretability.

Equation (4) modifies the original SCM proposal in two ways. First, Equation (4) excludes auxiliary covariates; we re-introduce them in Section 6. Second, Equation (4) penalizes the dispersion of the weights with hyperparameter $\zeta \geq 0$, following a suggestion in Abadie et al. (2015). Possible penalization functions include an elastic net penalty (Doudchenko and Imbens, 2017), an entropy penalty (Robbins et al., 2017), and a penalty based on a measure of pairwise distance (Abadie and L'Hour, 2018). The choice of penalty is less central when weights are constrained to be on the simplex, but becomes more important below when we relax this constraint (Doudchenko and Imbens, 2017).

We can view the SCM optimization problem in Equation (4) as an approximate balancing weights estimator (see, e.g. Zubizarreta, 2015; Hirshberg et al., 2019). As with all balancing estimators, a central question is what quantity to balance. Following the recent methodological literature (see Doudchenko and Imbens, 2017; Ferman and Pinto, 2018), Equation (4) directly op-

timizes for the pre-treatment fit, minimizing the (possibly weighted) imbalance of pre-treatment outcomes between the treated unit and the weighted control mean. In Section 5, we argue that this a natural quantity to target under both linearity and a latent factor model. Many choices are possible, however, and we can easily modify Equation (4) to balance other summary measures and functions of the lagged outcomes; see, for example, [Gobillon and Magnac \(2016\)](#), [Amjad et al. \(2018\)](#), and [Hazlett and Xu \(2018\)](#).

When the treated unit’s vector of lagged outcomes, \mathbf{X}_1 , is inside the convex hull of the control units’ lagged outcomes, \mathbf{X}_0 , the SCM weights in Equation (4) achieve perfect pre-treatment fit, and the resulting estimator has many attractive properties, including a bias bound established by [Abadie et al. \(2010\)](#). Due to the curse of dimensionality, however, achieving perfect (or nearly perfect) pre-treatment fit is not always feasible with weights constrained to be on the simplex (see [Ferman and Pinto, 2018](#)). When “the pre-treatment fit is poor or the number of pre-treatment periods is small,” [Abadie et al. \(2015\)](#) recommend against using SCM. Even if the pre-treatment fit is excellent, [Abadie et al. \(2010, 2015\)](#) propose extensive placebo checks to ensure that SCM weights do not overfit to noise. Thus, the conditional nature of the analysis is critical to deploying SCM, excluding many practical settings. Our proposal enables the use of (a modified) SCM approach in many of the cases where SCM alone is infeasible.

3 Augmented SCM

3.1 Overview

We now show how to modify the SCM approach to adjust for poor pre-treatment fit. Let \hat{m}_{iT} be an estimator for the post-treatment control potential outcome $Y_{iT}(0)$. The *Augmented SCM* (ASCM) estimator for $Y_{1T}(0)$ is:

$$\hat{Y}_{1T}^{\text{aug}}(0) = \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} Y_{iT} + \left(\hat{m}_{1T} - \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} \hat{m}_{iT} \right) \quad (5)$$

$$= \hat{m}_{1T} + \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} (Y_{iT} - \hat{m}_{iT}), \quad (6)$$

where weights $\hat{\gamma}_i^{\text{scm}}$ are the SCM weights defined above. Standard SCM is a special case, where \hat{m}_{iT} is a constant. We will largely focus on estimators that are a function of pre-treatment outcomes, $\hat{m}_{iT} \equiv \hat{m}(\mathbf{X}_i)$, where $\hat{m} : \mathbb{R}^{T_0} \rightarrow \mathbb{R}$.

Equations (5) and (6), while equivalent, highlight two distinct motivations for ASCM. Equation (5) directly corrects the SCM estimate, $\sum \hat{\gamma}_i^{\text{scm}} Y_{iT}$, by the imbalance in a particular function of the pre-treatment outcomes $\hat{m}(\cdot)$. Intuitively, since \hat{m} estimates the post-treatment outcome, we can view this as an estimate of the bias due to imbalance, analogous to bias correction for inexact

matching (Rubin, 1973; Abadie and Imbens, 2011). In this form, we can see that SCM and ASCM estimates will be similar if the estimated bias is small, as measured by imbalance in $\hat{m}(\cdot)$. In independent work, Abadie and L’Hour (2018) also consider a bias-corrected estimator of this form, though their paper focuses on reducing bias due to interpolation rather than extrapolation.

Equation (6), by contrast, is analogous to standard doubly robust estimation (Robins et al., 1994), which begins with the outcome model but then re-weights to balance residuals. This is also comparable in form to the generalized regression estimator in survey sampling (Cassel et al., 1976; Breidt and Opsomer, 2017), which has been adapted to the causal inference setting by, among others, Athey et al. (2018) and Hirshberg and Wager (2018). We discuss a connection to inverse propensity score weighting in Appendix E.

3.2 Choice of estimator

While this setup is general, the choice of estimator \hat{m} is important both for understanding the procedure’s properties and for practical performance. We give a brief overview of two special cases: (1) when \hat{m} is linear in the pre-treatment outcomes; and (2) when \hat{m} is linear in the comparison units. Ridge regression is an important example that is linear in both pre-treatment outcomes and comparison units; we explore this estimator further in Sections 4 and 5.

First, consider a model that is linear in pre-treatment outcomes, $\hat{m}(\mathbf{X}) = \hat{\eta}_0 + \hat{\boldsymbol{\eta}} \cdot \mathbf{X}$. The augmented estimator (5) is then:

$$\hat{Y}_{1T}^{\text{aug}}(0) = \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} Y_{iT} + \sum_{t=1}^{T_0} \hat{\eta}_t \left(X_{1t} - \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} X_{it} \right). \quad (7)$$

Pre-treatment periods that are more predictive of the post-treatment outcome will have larger (absolute) regression coefficients and so imbalance in these periods will lead to a larger adjustment. Thus, even if we do not *a priori* prioritize balance in any particular pre-treatment time periods (via the choice of \mathbf{V}_x), the linear model augmentation will adjust for the time periods that are empirically more predictive of the post-treatment outcome. As we show in Section 4, the ridge-regularized linear model is an important special case in which the resulting augmented estimator is itself a penalized synthetic control estimator. This allows for a more direct analysis of the role of bias correction.

Second, consider an outcome model that is a linear combination of comparison units, $\hat{m}(\mathbf{X}) = \sum_{W_i=0} \hat{\alpha}_i(\mathbf{X}) Y_{iT}$, for some weighting function $\hat{\alpha} : \mathbb{R}^{T_0} \rightarrow \mathbb{R}^{N_0}$. Examples include k -nearest neighbor matching and kernel weighting as well as other “vertical” regression approaches (Athey et al., 2017).

The augmented estimator (5) is itself a weighting estimator that adjusts the SCM weights:¹

$$\hat{Y}_{1T}^{\text{aug}}(0) = \sum_{W_i=0} \left(\hat{\gamma}_i^{\text{scm}} + \hat{\gamma}_i^{\text{adj}} \right) Y_{iT}, \quad \text{where} \quad \hat{\gamma}_i^{\text{adj}} \equiv \hat{\alpha}_i(\mathbf{X}_1) - \sum_{W_j=0} \hat{\gamma}_j^{\text{scm}} \hat{\alpha}_i(\mathbf{X}_j). \quad (8)$$

Here the adjustment term for unit i , $\hat{\gamma}_i^{\text{adj}}$, is the imbalance in a unit i -specific transformation of the lagged outcomes that depends on the weighting function $\alpha(\cdot)$. While $\hat{\gamma}^{\text{scm}}$ are constrained to be on the simplex, the form of $\hat{\gamma}^{\text{adj}}$ makes clear that the overall weights can be negative.

There are many special cases to consider, though we do not explore them in depth. One is the linear-in-lagged-outcomes model with equal coefficients, $\hat{\eta}_t = \frac{1}{T_0}$, which estimates a fixed-effects outcome model as $\hat{m}(\mathbf{X}_i) = \bar{X}_i$. The corresponding treatment effect estimate adjusts for imbalance in all pre-treatment time periods equally, and yields a weighted difference-in-differences estimator:

$$\hat{\tau}^{\text{de}} = (Y_{1T} - \bar{X}_1) - \left(\sum_{W_i=0} \hat{\gamma}_i (Y_{iT} - \bar{X}_i) \right) = \frac{1}{T_0} \sum_{t=1}^{T_0} \left[(Y_{1T} - X_{1t}) - \left(\sum_{W_i=0} \hat{\gamma}_i (Y_{iT} - X_{it}) \right) \right]. \quad (9)$$

An augmented estimator of this form has appeared as the *de-meaned* or *intercept shift SCM* (Doudchenko and Imbens, 2017; Ferman and Pinto, 2018).² See also Arkhangelsky et al. (2019), who extend this to weight across both units and time.

In Section 7.1 we conduct a simulation study to inspect the performance of a range of estimators, including other penalized linear models, such as the LASSO, flexible machine learning models, such as random forests, and panel data methods, such as fixed effects models and low-rank matrix completion methods (Xu, 2017; Athey et al., 2017).

4 Ridge ASCM improves pre-treatment fit while penalizing extrapolation

We now inspect the algorithmic and numerical properties for the special case where $\hat{m}(\mathbf{X}_i)$ is estimated via a ridge-regularized linear model, which we refer to as *Ridge Augmented SCM* (Ridge ASCM). With Ridge ASCM, the estimator for the post-treatment outcome is $\hat{m}(\mathbf{X}_i) = \hat{\eta}_0^{\text{ridge}} + \mathbf{X}'_i \hat{\boldsymbol{\eta}}^{\text{ridge}}$, where $\hat{\eta}_0^{\text{ridge}}$ and $\hat{\boldsymbol{\eta}}^{\text{ridge}}$ are the coefficients of a ridge regression of control post-treatment

¹We thank an anonymous reviewer for suggesting this presentation.

²In these proposals, the SCM weights balance the *residual* outcomes $X_{it} - \bar{X}_i$ rather than the raw outcomes X_{it} . We further consider balancing residuals in Section 6.

outcomes \mathbf{Y}_{0T} on centered pre-treatment outcomes \mathbf{X}_0 , with penalty hyper-parameter λ^{ridge} :³

$$\left\{ \hat{\eta}_0^{\text{ridge}}, \hat{\boldsymbol{\eta}}^{\text{ridge}} \right\} = \arg \min_{\eta_0, \boldsymbol{\eta}} \frac{1}{2} \sum_{W_i=0} (Y_i - (\eta_0 + X_i' \boldsymbol{\eta}))^2 + \lambda^{\text{ridge}} \|\boldsymbol{\eta}\|_2^2. \quad (10)$$

The Ridge Augmented SCM estimator is then:

$$\hat{Y}_{1T}^{\text{aug}}(0) = \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} Y_{iT} + \left(\mathbf{X}_1 - \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} \mathbf{X}_i \right) \cdot \hat{\boldsymbol{\eta}}^{\text{ridge}}. \quad (11)$$

We first show that Ridge ASCM is a linear weighting estimator as in Equation (8). Unlike augmenting with other linear weighting estimators, when augmenting with ridge regression the implied weights are themselves the solution to a penalized synthetic control problem, as in Equation (4). Using this representation, we show that when the treated unit lies outside the convex hull of the control units, Ridge ASCM improves the pre-treatment fit relative to SCM alone by allowing for negative weights and extrapolating away from the convex hull.

Allowing for negative weights is an important departure from the original SCM proposal, which constrains weights to be on the simplex. Ridge regression alone also allows for negative weights, and may have negative weights even when the treated unit is inside of the convex hull. In contrast, Ridge ASCM directly penalizes distance from the sparse, non-negative SCM weights, controlling the amount of extrapolation by the choice of λ^{ridge} , and only resorts to negative weights if the treated unit is outside of the convex hull.

4.1 Ridge ASCM as a penalized SCM estimator

We now express both Ridge ASCM and ridge regression alone as special cases of the penalized SCM problem in Equation (4). The Ridge ASCM estimate of the counterfactual is the solution to (4), replacing the simplex constraint with a penalty $f(\gamma_i) = (\gamma_i - \hat{\gamma}_i^{\text{scm}})^2$ that penalizes *deviations from the SCM weights*.

Lemma 1. The ridge-augmented SCM estimator (7) is:

$$\hat{Y}_{1T}^{\text{aug}}(0) = \sum_{W_i=0} \hat{\gamma}_i^{\text{aug}} Y_{iT}, \quad (12)$$

where

$$\hat{\gamma}_i^{\text{aug}} = \hat{\gamma}_i^{\text{scm}} + (\mathbf{X}_1 - \mathbf{X}'_0 \hat{\boldsymbol{\gamma}}^{\text{scm}})' (\mathbf{X}'_0 \mathbf{X}_0 + \lambda^{\text{ridge}} \mathbf{I}_{T_0})^{-1} \mathbf{X}_i. \quad (13)$$

³Similar to the synthetic controls problem, we can regularize time periods differently with a generalized ridge penalty $\boldsymbol{\eta}' \boldsymbol{\Lambda} \boldsymbol{\eta}$ using an importance matrix $\boldsymbol{\Lambda}$. Following the typical case with diagonal elements, the generalized ridge penalty reduces to separate regularization on each time period.

Moreover, the Ridge ASCM weights $\hat{\gamma}^{\text{aug}}$ are the solution to

$$\min_{\gamma \text{ s.t. } \sum_i \gamma_i = 1} \frac{1}{2\lambda^{\text{ridge}}} \|\mathbf{X}_1 - \mathbf{X}'_0 \gamma\|_2^2 + \frac{1}{2} \|\gamma - \hat{\gamma}^{\text{scm}}\|_2^2. \quad (14)$$

When the treated unit is in the convex hull of the control units — so the SCM weights exactly balance the lagged outcomes — the Ridge ASCM and SCM weights are identical. When SCM weights do not achieve exact balance, the Ridge ASCM solution will use negative weights to extrapolate from the convex hull of the control units. The amount of extrapolation is determined both by the amount of imbalance and by the hyperparameter λ^{ridge} . When SCM yields good pre-treatment fit or when λ^{ridge} is large, the adjustment term will be small and $\hat{\gamma}^{\text{aug}}$ will remain close to the SCM weights.

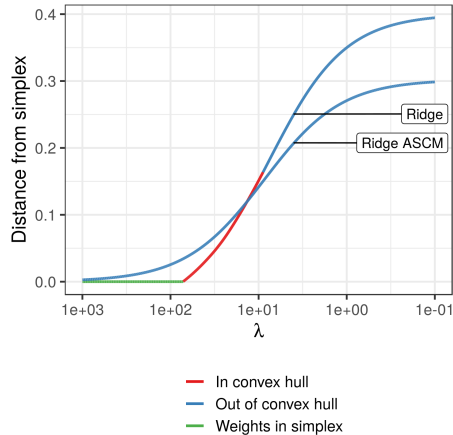
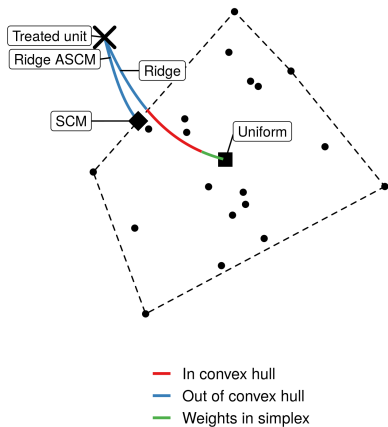
We can similarly characterize ridge regression alone as a solution to a penalized SCM problem where the penalty term, $f(\gamma_i) = \left(\gamma_i - \frac{1}{N_0}\right)^2$, penalizes the variance of the weights. Other penalized linear models, such as the LASSO or elastic net, do not have this same representation as a penalized SCM estimator.

Lemma 2. The ridge regression estimator $\hat{Y}_{1T}^{\text{ridge}}(0) \equiv \hat{\eta}_0^{\text{ridge}} + \mathbf{X}_1 \cdot \hat{\boldsymbol{\eta}}^{\text{ridge}}$ can be written as $\hat{Y}_{1T}^{\text{ridge}}(0) = \sum_{W_i=0} \hat{\gamma}_i^{\text{ridge}} Y_{iT}$, where the ridge weights $\hat{\boldsymbol{\gamma}}^{\text{ridge}}$ are the solution to:

$$\min_{\gamma \mid \sum_i \gamma_i = 1} \frac{1}{2\lambda^{\text{ridge}}} \|\mathbf{X}_1 - \mathbf{X}'_0 \gamma\|_2^2 + \frac{1}{2} \left\| \gamma - \frac{1}{N_0} \right\|_2^2. \quad (15)$$

For ridge regression alone, the hyperparameter λ^{ridge} controls the variance of the weights rather than the degree of extrapolation from the simplex. Thus, in order to reduce variance, the ridge regression weights might still be negative even if the treated unit is inside of the convex hull and SCM achieves perfect fit.

Figure 1 visualizes these results in two dimensions. Figure 1a shows the distribution of control units and the treated unit outside the convex hull, along with the weighted average of control units using the ridge and Ridge ASCM weights, varying λ^{ridge} . We see that ridge regression alone begins (for large λ^{ridge}) at the center of the control units, while Ridge ASCM begins at the SCM solution; both move smoothly towards an exact fit solution as λ^{ridge} is reduced. Figure 1b shows the distance of the ridge and Ridge ASCM weights from the simplex as λ^{ridge} varies. Ridge ASCM weights begin at the SCM solution, which is on the boundary of the simplex, then extrapolate outside the convex hull. In contrast, ridge regression weights begin in the center of the simplex (i.e., uniform weights), but then leave the simplex (i.e., some negative weights) before the corresponding weighted average is outside of the convex hull. Over this range, marked as red in Figure 1a, it is possible to achieve the same level of balance with non-negative weights, but ridge regression uses negative weights in order to reduce the variance. Eventually, as $\lambda^{\text{ridge}} \rightarrow 0$, both ridge and Ridge ASCM use negative weights to achieve perfect balance; the weight vectors are different, however, with the Ridge ASCM



(a) Treated and control units with the convex hull marked as a dashed line. Ridge and Ridge ASCM estimates in solid. (b) Distance of ridge and Ridge ASCM weights from the simplex.

Figure 1: Ridge ASCM vs. ridge regression alone for a two-dimensional example with the treated unit outside of the convex hull of the control units. Results shown varying λ^{ridge} from 10^3 to 10^{-1} . Green denotes that the weights are inside the simplex, red that the weights are outside the simplex but the weighted average is inside of the convex hull, and blue that the weighted average is outside the convex hull.

vector closer to the simplex.

As this example makes clear, both ridge and Ridge ASCM extrapolate from the support of the data to improve pre-treatment fit relative to SCM alone. Ridge ASCM, however, is equivalent to SCM weights when SCM achieves exact pre-treatment fit. When achieving excellent pre-treatment fit with SCM is possible, [Abadie et al. \(2015\)](#) argue that we should prefer SCM weights over possibly negative weights: a slight balance improvement is not worth the extrapolation and the loss of interpretability. In this case, the Ridge ASCM weights will be close to the simplex, while the ridge regression weights may be quite far away. When this is not possible, however, and SCM has poor fit, some degree of extrapolation is critical; Ridge ASCM allows the researcher to directly penalize the amount of extrapolation in these cases.⁴

4.2 Ridge ASCM improves pre-treatment fit relative to SCM alone

Just as the hyper-parameter λ^{ridge} parameterizes the level of extrapolation, it also parameterizes the level of improvement in pre-treatment fit over the SCM solution. Because we are removing the non-negativity constraint and allowing for extrapolation outside of the convex hull, the pre-treatment fit from Ridge ASCM will be at least as good as the pre-treatment fit from SCM alone,

⁴See [King and Zeng \(2006\)](#) for a discussion of extrapolation in constructing counterfactuals. As they note: “If we learn that a counterfactual question involves extrapolation, we still might wish to proceed if the question is sufficiently important, but we would be aware of how much more model dependent our answers would be.”

i.e., $\|\mathbf{X}_1 - \mathbf{X}'_0 \hat{\gamma}^{\text{aug}}\|_2 \leq \|\mathbf{X}_1 - \mathbf{X}'_0 \hat{\gamma}^{\text{scm}}\|_2$. We can exactly characterize the pre-treatment fit of Ridge ASCM using the singular value decomposition of the matrix of control outcomes.

Lemma 3. Let $\frac{1}{\sqrt{N_0}}\mathbf{X}_0 = \mathbf{U}\mathbf{D}\mathbf{V}'$ be the singular value decomposition of the matrix of control pre-intervention outcomes, where m is the rank of \mathbf{X}_0 , $\mathbf{U} \in \mathbb{R}^{N_0 \times m}$, $\mathbf{V} \in \mathbb{R}^{T_0 \times m}$, and $\mathbf{D} = \text{diag}(d_1, \dots, d_m) \in \mathbb{R}^{m \times m}$ is the diagonal matrix of singular values, where d_1 and d_m are the largest and smallest singular values, respectively. Furthermore, let $\tilde{\mathbf{X}}_i = \mathbf{V}'\mathbf{X}_i$ be the rotation of \mathbf{X}_i along the singular vectors of \mathbf{X}_0 . Then $\hat{\gamma}^{\text{aug}}$, the Ridge ASCM weights with hyper-parameter $\lambda^{\text{ridge}} = \lambda N_0$ satisfy

$$\|\mathbf{X}_1 - \mathbf{X}'_0 \hat{\gamma}^{\text{aug}}\|_2 = \lambda \left\| (\mathbf{D} + \lambda \mathbf{I})^{-1} (\tilde{\mathbf{X}}_1 - \tilde{\mathbf{X}}'_0 \hat{\gamma}^{\text{scm}}) \right\|_2 \leq \frac{\lambda}{d_m^2 + \lambda} \|\mathbf{X}_1 - \mathbf{X}'_0 \hat{\gamma}^{\text{scm}}\|_2, \quad (16)$$

and the weights from ridge regression alone $\hat{\gamma}^{\text{ridge}}$ satisfy

$$\|\mathbf{X}_1 - \mathbf{X}'_0 \hat{\gamma}^{\text{ridge}}\|_2 = \lambda \left\| (\mathbf{D} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}_1 \right\|_2 \leq \frac{\lambda}{d_m^2 + \lambda} \|\mathbf{X}_1\|_2. \quad (17)$$

From Equation (16), we see that the pre-treatment imbalance for Ridge ASCM weights is smaller than that of SCM weights by at least a factor of $\frac{\lambda}{d_m^2 + \lambda}$. Thus, Ridge ASCM will achieve strictly better pre-treatment fit than SCM alone, except in corner cases where pre-treatment fit will be equal. For example, these will be equal when the pre-treatment SCM residual $\mathbf{X}_1 - \mathbf{X}'_0 \hat{\gamma}^{\text{scm}}$ is orthogonal to the lagged outcomes of the control units \mathbf{X}_0 . From Equation (17), we see that the relationship for pre-treatment imbalance fit between SCM and ridge regression is less clear; ridge regression penalizes deviations from uniformity, rather than deviations from SCM weights.

5 Error under restrictions on the data generating processes

We now relate Ridge ASCM's improved pre-treatment fit to improved estimation error under several data generating processes. We first consider the case where the post-treatment outcome is linear in the pre-treatment outcomes, and then extend to the canonical linear factor model considered by [Abadie et al. \(2010\)](#). In the appendix, we also consider a more general formulation, including when the outcome model is approximately linear, with Lipschitz deviations from linearity.

Under a linear model, improving pre-treatment fit directly reduces bias, and the Ridge ASCM penalty term negotiates a bias-variance trade-off. Under a latent factor model, improving pre-treatment fit again reduces bias, though there is now a risk of over-fitting. The penalty term also directly parameterizes this trade-off. Thus, choosing the hyper-parameter λ^{ridge} is important in practice. In Section 5.3, we describe a cross-validation hyper-parameter selection procedure that builds on the in-time placebo check in [Abadie et al. \(2015\)](#). In Section 7.1 we explore these trade-offs through simulations, finding substantial gains to augmentation via ridge regression in a variety

of settings.

5.1 Error under linearity

We first illustrate the key balancing idea in the simple case where the post-treatment outcome is a linear combination of lagged outcomes plus additive noise:

$$Y_{iT}(0) = \sum_{t=1}^{T_0} \beta_t \mathbf{X}_{it} + \varepsilon_{iT}. \quad (18)$$

A special case of this setup is an auto-regressive process of order $K \leq T_0$. Here we consider the pre-treatment outcomes \mathbf{X}_i fixed and so the randomness in the post-treatment outcome $Y_{iT}(0)$ is due to the noise term ε_{iT} . As in Section 2, we assume that the N units' noise terms are independent, mean zero sub-Gaussian random variables with scale parameter σ that are uncorrelated with treatment assignment W_i .

We consider a generic weighting estimator with weights $\hat{\gamma}$ that are independent of the post-treatment outcomes Y_{1T}, \dots, Y_{NT} ; both SCM and Ridge ASCM take this form. Under linearity, the difference between the counterfactual outcome $Y_{1T}(0)$ and the weighting estimator $\hat{Y}_{1T}(0)$ decomposes into: (1) systemic error due to imbalance in the lagged outcomes \mathbf{X} , and (2) idiosyncratic error due to the noise in the post-treatment period:

$$Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i Y_{iT} = \underbrace{\boldsymbol{\beta} \cdot \left(\mathbf{X}_1 - \sum_{W_i=0} \mathbf{X}_i \right)}_{\text{imbalance in } \mathbf{X}} + \underbrace{\varepsilon_{1T} - \sum_{W_i=0} \hat{\gamma}_i \varepsilon_{iT}}_{\text{post-treatment noise}}. \quad (19)$$

With this setup, a weighting estimator that exactly balances the lagged outcomes \mathbf{X} will eliminate all systematic error. Furthermore, if $\boldsymbol{\beta}$ is sparse, then it suffices to balance only the lagged outcomes with non-zero coefficients; for example, under an AR(K) process, $(\beta_1, \dots, \beta_{T_0-K-1}) = 0$, it is sufficient to balance only the first K lags.

If the weighting estimator does not perfectly balance the pre-treatment outcomes \mathbf{X} , there will be a systematic component of the error, with the magnitude depending on the imbalance. Below we construct a finite sample error bound for Ridge ASCM (and for SCM, the special case with $\lambda^{\text{ridge}} = \infty$). This bound on the estimation error holds with high probability over the noise in the post-treatment period $\boldsymbol{\varepsilon}_T$.

Proposition 1. Under the linear model (18) with independent sub-Gaussian noise with scale parameter σ , for any $\delta > 0$, the Ridge ASCM weights with hyperparameter $\lambda^{\text{ridge}} = \lambda N_0$ satisfy

the bound

$$\left| Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i^{\text{aug}} Y_{iT} \right| \leq \|\beta\|_2 \underbrace{\left\| \text{diag} \left(\frac{\lambda}{d_j^2 + \lambda} \right) (\widetilde{\mathbf{X}}_1 - \widetilde{\mathbf{X}}_0' \hat{\gamma}^{\text{scm}}) \right\|_2}_{\text{imbalance in } \mathbf{X}} + \underbrace{\delta\sigma (1 + \|\hat{\gamma}^{\text{aug}}\|_2)}_{\text{post-treatment noise}}, \quad (20)$$

with probability at least $1 - 2e^{-\frac{\delta^2}{2}}$, where $\widetilde{\mathbf{X}}_i = \mathbf{V}' \mathbf{X}_i$ is the rotation of \mathbf{X}_i along the singular vectors of \mathbf{X}_0 , as above.

Proposition 1 shows the finite sample error of Ridge ASCM weights is controlled by the imbalance in the lagged outcomes and the L^2 norm of the weights; Lemma A.2 in the Appendix gives a deterministic bound for $\|\hat{\gamma}^{\text{aug}}\|_2$. See Athey et al. (2018) for analogous results on balancing weights in high dimensional cross-sectional settings.

In the special case that SCM weights have perfect pre-treatment fit, ASCM and SCM weights will be equivalent, and the estimation error will only be due to variance of the weights and post-treatment noise. When SCM weights do not achieve perfect pre-treatment fit, Ridge ASCM with finite $\lambda < \infty$ extrapolates outside the convex hull, improving pre-treatment fit and thus reducing bias. This is subject to the usual bias-variance trade-off: The second term in (20) is increasing in the L^2 norm of the weights, which will generally be larger for ASCM than for SCM. The hyperparameter λ directly negotiates this trade off.

5.2 Error under a latent factor model

We now consider the case where control potential outcomes are generated according to a linear factor model. Similar to the linear case above, when pre-treatment fit is poor, improving the pre-treatment fit by extrapolating away from the convex hull reduces the bias, with a possible bias-variance tradeoff. Unlike in the linear setting, however, there is additional possible error due to balancing noisy pre-treatment outcomes rather than the underlying factors. This raises the risk of over-fitting, though it can be reduced by appropriate choice of the extrapolation penalty. Specifically, in settings where the noise is high, ASCM can increase error relative to SCM alone. SCM limits over-fitting in this high-noise case by constraining weights to the simplex, although SCM alone is likely to perform poorly here regardless.

Following the setup in Abadie et al. (2010), we assume that there are J unknown, latent time-varying factors $\boldsymbol{\mu}_t = \{\mu_{jt}\} \in \mathbb{R}^T$, $j = 1, \dots, J$, with $\max_{jt} |\mu_{jt}| \leq M$, where J will typically be small relative to N and T_0 . In addition, each unit has a vector of unknown factor loadings $\phi_i \in \mathbb{R}^J$. We consider both the time-varying factors $\boldsymbol{\mu}_t$ and the unit-varying factor loadings ϕ_i to be non-random quantities and will consider a fixed N and T_0 . In this setup, the control potential outcomes are

weighted averages of these factors plus additive noise:

$$Y_{it}(0) = \sum_{j=1}^J \phi_{ij} \mu_{jt} + \varepsilon_{it} = \boldsymbol{\phi}_i \cdot \boldsymbol{\mu}_t + \varepsilon_{it}, \quad (21)$$

where again the randomness in $Y_{it}(0)$ is only due to the noise term ε_{it} . We assume that ε_{it} are independent, mean zero sub-Gaussian random variables with scale parameter σ that are uncorrelated with treatment assignment W_i . Slightly abusing notation, we collect the pre-intervention factors into a matrix $\boldsymbol{\mu} \in \mathbb{R}^{T_0 \times J}$, where the t^{th} row of $\boldsymbol{\mu}$ contains the factor values at time t , $\boldsymbol{\mu}'_t$. Following Bai (2009) and Xu (2017), we further assume that the factors are orthogonal and normalized, i.e., that $\frac{1}{T_0} \boldsymbol{\mu}' \boldsymbol{\mu} = \mathbf{I}_J$.

Under this model, the finite-sample error of a weighting estimator depends on the imbalance in the latent $\boldsymbol{\phi}$ and a noise term due to the noise at time T :

$$Y_{1T}(0) - \hat{Y}_{1T}(0) = Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i Y_{iT} = \underbrace{\left(\boldsymbol{\phi}_1 - \sum_{W_i=0} \hat{\gamma}_i \boldsymbol{\phi}_i \right) \cdot \boldsymbol{\mu}_T}_{\text{imbalance in } \boldsymbol{\phi}} + \underbrace{\varepsilon_{1T} - \sum_{W_i=0} \hat{\gamma}_i \varepsilon_{it}}_{\text{noise}}. \quad (22)$$

Balancing the observed pre-treatment outcomes \mathbf{X} will not necessarily balance the latent factor loadings $\boldsymbol{\phi}$. Following Abadie et al. (2010), we show in the appendix that, under Equation (21), we can decompose the imbalance term as:

$$\left(\boldsymbol{\phi}_1 - \sum_{W_i=0} \gamma_i \boldsymbol{\phi}_i \right) \cdot \boldsymbol{\mu}_T = \frac{1}{T_0} \boldsymbol{\mu}' \underbrace{\left(\mathbf{X}_1 - \sum_{W_i=0} \gamma_i \mathbf{X}_i \right)}_{\text{imbalance in } \mathbf{X}} \cdot \boldsymbol{\mu}_T - \frac{1}{T_0} \boldsymbol{\mu}' \underbrace{\left(\boldsymbol{\varepsilon}_{1(1:T_0)} - \sum_{W_i=0} \gamma_i \boldsymbol{\varepsilon}_{i(1:T_0)} \right)}_{\text{approximation error}} \cdot \boldsymbol{\mu}_T, \quad (23)$$

where $\boldsymbol{\varepsilon}_{i(1:T_0)} = (\varepsilon_{i1}, \dots, \varepsilon_{iT_0})$ is the vector of pre-treatment noise terms for unit i . The first term is the imbalance of observed lagged outcomes and the second term is an approximation error arising from the latent factor structure. With $\sigma = 0$ and J small, the approximation error is zero, and it is possible to express $Y_{iT}(0)$ as a linear combination of $Y_{it}(0)$, $t = 1, \dots, T_0 - 1$, recovering the linear model (18). However, with $\sigma > 0$ we cannot write the period- T outcome as a linear combination of earlier outcomes plus independent, additive error.

With this setup, we can bound the finite-sample error in Equation (22) for Ridge ASCM weights (and for SCM weights as a special case). This bound is with high probability over the noise in all time periods ε_{it} , and accounts for the noise in the pre- and post-treatment outcomes separately.

Theorem 1. Under the linear factor model (21) with independent sub-Gaussian noise, for any

$\delta > 0$, the Ridge ASCM weights with hyperparameter $\lambda^{\text{ridge}} = \lambda N_0$ satisfy the bound

$$\begin{aligned}
\left| Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i^{\text{aug}} Y_{1T}(0) \right| &\leq \frac{JM^2}{\sqrt{T_0}} \left(\underbrace{\left\| \text{diag} \left(\frac{\lambda}{d_j^2 + \lambda} \right) (\widetilde{\mathbf{X}}_1 - \widetilde{\mathbf{X}}'_0 \hat{\gamma}^{\text{scm}}) \right\|_2}_{\text{imbalance in } \mathbf{X}} + \right. \\
&\quad \underbrace{4(1 + \delta) \left\| \text{diag} \left(\frac{d_j \sigma}{d_j^2 + \lambda} \right) (\widetilde{\mathbf{X}}_1 - \widetilde{\mathbf{X}}'_0 \hat{\gamma}^{\text{scm}}) \right\|_2}_{\text{excess approximation error}} + \\
&\quad \left. \underbrace{2 \left(\sqrt{\log 2N_0} + \frac{\delta}{2} \right)}_{\text{SCM approximation error}} \right) + \underbrace{\delta \sigma (1 + \|\hat{\gamma}^{\text{aug}}\|_2)}_{\text{post-treatment noise}}
\end{aligned} \tag{24}$$

with probability at least $1 - 6e^{-\frac{\delta^2}{2}} - e^{-2(\log 2 + N_0 \log 5)\delta^2}$.

Theorem 1 shows that, relative to the linear case in Proposition 1, there is an additional source of error under a latent factor model: approximation error due to balancing lagged outcomes rather than balancing underlying factors. In particular, it is now possible that a control unit only receives a large weight because of idiosyncratic noise, rather than because of similarity in the underlying factors. See Arkhangelsky et al. (2019) and Ferman (2019) for asymptotic analogues of this finite sample bound.

In the special case where SCM achieves perfect pre-treatment fit, considered by Abadie et al. (2010), the ASCM and SCM weights are equivalent and the error is only due to post-treatment noise and the approximation error. The bound in Theorem 1 accounts for the worst case scenario where the control unit with the largest weight is only similar to the treated unit due to idiosyncratic noise. The approximation error, and thus the bias, converges to zero in probability as $T_0 \rightarrow \infty$ under suitable conditions on the factor loadings $\boldsymbol{\mu}_t$ (see also Ferman and Pinto, 2018). Intuitively, as we observe more X_{it} — and can exactly balance each one — we are better able to match on the index $\boldsymbol{\phi}_i \cdot \boldsymbol{\mu}_t$ and, as a result, on the underlying factor loadings.⁵

Without exact balance, Theorem 1 shows that a long pre-period may not be enough to control the error due to imbalance. This emphasizes the critical role that conditioning on excellent pre-treatment fit plays in ensuring that SCM alone yields reasonable estimates of the counterfactual.

When perfect pre-treatment fit is not feasible with SCM, Ridge ASCM with $\lambda < \infty$ will extrapolate outside the convex hull, where the hyper-parameter λ controls the amount of extrapolation. As in the linear case, extrapolation reduces the error due to imbalance in lagged outcomes. Unlike

⁵We show in the supplementary material that with dependent errors the probability of the worst-case error additionally scales with the maximum eigenvalue of the covariance matrix. Dependence leads to a more complicated error structure overall; we leave a thorough analysis of this to future work.

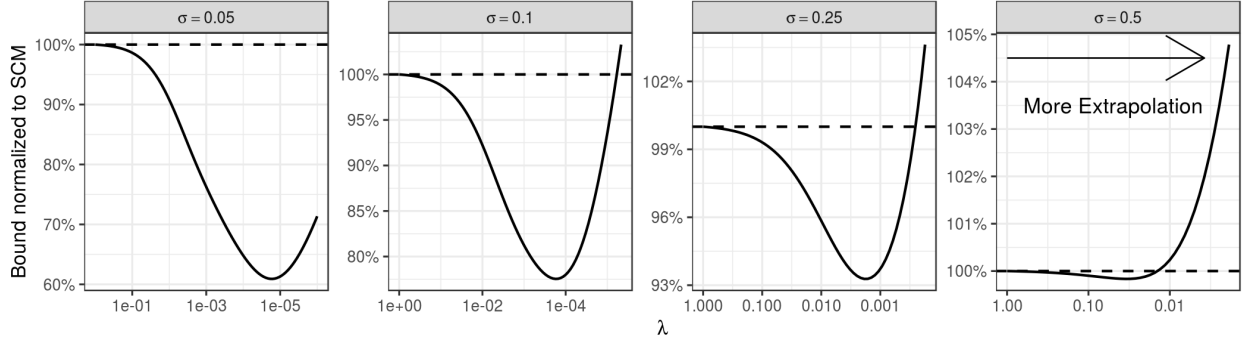


Figure 2: Sketch of the error due to imbalance and approximation error (24) for the linear factor model; the standard deviation of the treated unit’s pre-treatment outcomes is normalized to one. We fit SCM weights on the empirical example in Section 7 and compute the vector of pre-treatment fit. Each line shows the sum of the error due to imbalance in \mathbf{X} , excess approximation error, and SCM approximation error in Theorem 1 (with $\delta = 0$) for different values of σ . These are normalized so that the SCM solution (with λ large) equals 100%; values below 100% show improvement over the unadjusted weights for a given λ .

in the linear case, however, this could possibly lead to over-fitting to the noisy lagged outcomes. While some extrapolation can be helpful on net, the optimal amount will depend on the synthetic control fit and the amount of noise. Figure 2 illustrates this using SCM weights from the empirical example we discuss in Section 7, where pre-treatment fit is relatively poor. For each value of σ , the figure plots the sum of the imbalance, SCM approximation error, and excess approximation error terms in the bound in Theorem 1. At each noise level, a small amount of extrapolation leads to a smaller error bound, but as λ shrinks there is a point where further extrapolation leads to overfitting and eventually to a worse error bound than without extrapolation. The risk of overfitting is greater when the noise is particularly large, though even here a sufficiently regularized ASCM estimate has lower error bound than SCM alone (represented as the $\lambda \rightarrow \infty$ bound). When noise is less extreme, the benefits of augmentation are larger and the optimal amount of regularization shrinks.

It is worth noting that Theorem 1 gives a worst-case bound. In Section 7.1 we inspect the typical performance of the Ridge ASCM estimator via extensive simulation studies and find that gains to pre-treatment fit through augmentation outweigh increased approximation error in a range of practical settings, including when noise is very large.

5.3 Hyper-parameter selection

We propose a cross-validation approach for selecting λ inspired by the in-time placebo check proposed by Abadie et al. (2015). Let $\hat{Y}_{1k}^{(-t)} = \sum_{W_i=0} \hat{\gamma}_{i(-t)}^{\text{aug}} Y_{ik}$ be the estimate of Y_{1k} where time period t is excluded from fitting the estimator in (13). Abadie et al. (2015) propose to compare

the difference $Y_{1t} - \hat{Y}_{1t}^{(-t)}$ for some $t \leq T_0$ as a placebo check. We can extend this idea to compute the leave-one-out cross validation MSE over time periods:

$$CV(\lambda) = \sum_{t=1}^{T_0} \left(Y_{1t} - \hat{Y}_{1t}^{(-t)} \right)^2. \quad (25)$$

We can then choose λ to minimize $CV(\lambda)$ or follow a more conservative approach such as the “one-standard-error” rule (Hastie et al., 2009). This proposal is similar to the leave-one-out cross validation proposed by Doudchenko and Imbens (2017), who select hyperparameters by holding out control units and minimizing the MSE of the control units in the post-treatment time T . Finally, only excluding time period t might be inappropriate for some outcome models, e.g. the linear model in Section 5.1. In these settings we can extend the procedure to exclude all time periods $\geq t$ when estimating $\hat{\gamma}_{(-t)}^{\text{aug}}$. For related proposals, see Kellogg et al. (2020) and Bilinski and Hatfield (2020).

6 Auxiliary covariates

Thus far, we have focused exclusively on lagged outcomes as predictors. We now consider the case where there are also a small number of auxiliary covariates $\mathbf{Z}_i \in \mathbb{R}^K$ for unit i . These auxiliary covariates may include summaries of lagged outcomes or time-varying covariates such as the pre-treatment mean \bar{X}_i . Let $\mathbf{Z}_0 \in \mathbb{R}^{N_0 \times K}$ denote the matrix of covariates for the donor units. We assume that the covariates are centered, $\bar{\mathbf{Z}}_0 = \mathbf{0}$.

These auxiliary covariates can be incorporated both into the balance objective for SCM and into the outcome model used for augmentation in ASCM. For example, we can extend SCM to choose weights to solve

$$\min_{\gamma \in \Delta^{N_0}} \theta_x \|\mathbf{X}_1 - \mathbf{X}'_0 \gamma\|_2^2 + \theta_z \|\mathbf{Z}_1 - \mathbf{Z}_0 \gamma\|_2^2 + \zeta \sum_{W_i=0} f(\gamma_i), \quad (26)$$

where Δ^{N_0} is the N_0 -simplex. Similarly, we can augment these weights with an outcome model $\hat{m}(\mathbf{X}_i, \mathbf{Z}_i)$ that is a function of both the lagged outcomes and auxiliary covariates. For example, we can extend Ridge ASCM to choose $\hat{m}(\mathbf{X}, \mathbf{Z}) = \hat{\eta}_0 + \mathbf{X}' \hat{\boldsymbol{\eta}}_x + \mathbf{Z}' \hat{\boldsymbol{\eta}}_z$ and fit via ridge regression:

$$\min_{\eta_0, \boldsymbol{\eta}_x, \boldsymbol{\eta}_z} \frac{1}{2} \sum_{W_i=0} (Y_i - (\eta_0 + \mathbf{X}'_i \boldsymbol{\eta}_x + \mathbf{Z}'_i \boldsymbol{\eta}_z))^2 + \lambda_x \|\boldsymbol{\eta}_x\|_2^2 + \lambda_z \|\boldsymbol{\eta}_z\|_2^2. \quad (27)$$

Both this SCM criterion and augmentation estimator incorporate user-specified weights that determine the importance of balancing each set of covariates (Equation 26) or the amount of regularization for each set of coefficients (Equation 27). There are many potential choices for these weights. We discuss two, appropriate to different settings depending on the number of auxiliary covariates.

A sensible default when the dimension of the auxiliary covariates is moderate is to incorporate

the lagged outcomes \mathbf{X} and the auxiliary covariates \mathbf{Z} equally in Equations (26) and (27), setting $\theta_x = \theta_z = 1$ and $\lambda_x = \lambda_z = \lambda^{\text{ridge}}$, after standardizing auxiliary covariates and lagged outcomes to have the same standard deviation. With this setup the numerical and algorithmic results in Section 4 apply for the combined vector of lagged outcomes and auxiliary covariates, $(\mathbf{X}_i, \mathbf{Z}_i) \in \mathbb{R}^{T_0+K}$. In particular, Ridge ASCM is again a penalized SCM estimator that adjusts the synthetic control weights that solve optimization problem (26) to achieve better balance by extrapolating outside of the convex hull.

An alternative approach when the dimension of the auxiliary covariates is small relative to N (i.e., $K \ll N$) is to fit a regression model that regularizes the lagged outcome coefficients $\boldsymbol{\eta}_x$ but does *not* regularize the auxiliary covariate coefficients $\boldsymbol{\eta}_z$ (i.e., set $\lambda_z = 0$). Lemma 4 below writes the resulting augmented estimator as its corresponding penalized SCM optimization problem, with weights that perfectly balance the auxiliary covariates. This has two key implications. First, since the auxiliary covariates \mathbf{Z} are exactly balanced regardless of the balance that the SCM weights achieve alone, we can exclude them from the optimization problem (26). Second, as we show below, the pre-treatment fit on the lagged outcomes depends on how well the SCM weights balance the residualized lagged outcomes $\check{\mathbf{X}}$. This suggests modifying Equation (26) to balance $\check{\mathbf{X}}$ rather than the lagged outcomes \mathbf{X} , which leads to the two-step procedure: (1) residualize the pre- and post-treatment outcomes on the auxiliary covariates \mathbf{Z} ; and (2) estimate Ridge ASCM on the residualized outcomes. This two-step procedure follows from a related proposal in [Doudchenko and Imbens \(2017\)](#).

Lemma 4. Let $\hat{\boldsymbol{\eta}}_x$ and $\hat{\boldsymbol{\eta}}_z$ be the solutions to (27) with $\lambda_x = \lambda^{\text{ridge}}$ and $\lambda_z = 0$. For any weight vector $\hat{\boldsymbol{\gamma}}$ that sums to one, the ASCM estimator from Equation (6) with $\hat{m}(\mathbf{X}_i, \mathbf{Z}_i) = \mathbf{X}_i' \hat{\boldsymbol{\eta}}_x + \mathbf{Z}_i' \hat{\boldsymbol{\eta}}_z$ is

$$\sum_{W_i=0} \hat{\gamma}_i Y_{iT} + \left(\mathbf{X}_1 - \sum_{W_i=0} \hat{\gamma}_i \mathbf{X}_i \right)' \hat{\boldsymbol{\eta}}_x + \left(\mathbf{Z}_1 - \sum_{W_i=0} \hat{\gamma}_i \mathbf{Z}_i \right)' \hat{\boldsymbol{\eta}}_z = \sum_{W_i=0} \hat{\gamma}_i^{\text{cov}} Y_{iT}, \quad (28)$$

where the weights $\hat{\boldsymbol{\gamma}}^{\text{cov}}$ are

$$\hat{\gamma}_i^{\text{cov}} = \hat{\gamma}_i + (\check{\mathbf{X}}_1 - \check{\mathbf{X}}_0)' (\check{\mathbf{X}}_0' \check{\mathbf{X}}_0 + \lambda^{\text{ridge}} \mathbf{I}_{T_0})^{-1} \check{\mathbf{X}}_i + (\mathbf{Z}_1 - \mathbf{Z}_0' \hat{\boldsymbol{\gamma}})' (\mathbf{Z}_0' \mathbf{Z}_0)^{-1} \mathbf{Z}_i, \quad (29)$$

and $\check{\mathbf{X}}_i$ is the residual components of a regression of pre-treatment outcomes on the control auxiliary covariates:

$$\check{\mathbf{X}}_i = \mathbf{X}_i - \mathbf{Z}_i' (\mathbf{Z}_0' \mathbf{Z}_0)^{-1} \mathbf{Z}_0' \mathbf{X}_0. \quad (30)$$

These weights exactly balance the auxiliary covariates, $\mathbf{Z}_1 - \mathbf{Z}_0' \hat{\boldsymbol{\gamma}}^{\text{cov}} = 0$; the imbalance in the lagged outcomes is

$$\| \mathbf{X}_1 - \mathbf{X}_0' \hat{\boldsymbol{\gamma}}^{\text{cov}} \|_2 \leq \left(\frac{\lambda^{\text{ridge}}}{\lambda^{\text{ridge}} + N_0 \check{d}_T^2} \right) \| \check{\mathbf{X}}_1 - \check{\mathbf{X}}_0' \hat{\boldsymbol{\gamma}} \|_2, \quad (31)$$

where \check{d}_r is the minimal singular value of $\check{\mathbf{X}}_0$.

Comparing to the numerical results in Section 4, Lemma 4 shows that the two-step approach penalizes extrapolation from the convex hull *in the residualized space* $\check{\mathbf{X}}$, rather than in the lagged outcomes themselves. In essence, by residualizing out the auxiliary covariates \mathbf{Z} the two-step approach allows for a possibly large amount of extrapolation in the auxiliary covariates, while carefully penalizing extrapolation in the part of the lagged outcomes that is orthogonal to the covariates.

In the Appendix, we consider the performance of this estimator when the outcomes follow a linear factor model with covariates for the special case where $\lambda^{\text{ridge}} \rightarrow \infty$ and the weights $\hat{\gamma}^{\text{cov}}$ do not extrapolate from the convex hull after residualization. We find that, if K is small relative to N_0 , exactly balancing a small number of auxiliary covariates and targeting imbalance in the residuals $\check{\mathbf{X}}$ can lead to decreased error due to pre-treatment fit, with only a small increase in approximation error and error due to post-treatment noise. However, with larger numbers of auxiliary covariates, the approach that incorporates auxiliary covariates in parallel to lagged outcomes would be more appropriate.

7 Simulations and empirical illustrations

We now turn to simulations and an empirical illustration. First, we conduct extensive simulation studies to assess the performance of different methods, finding substantial gains from ASCM. We then use our approach to examine the effect of an aggressive tax cut on economic output in Kansas in 2012.

7.1 Calibrated simulation studies

We now present simulation studies calibrated to our empirical illustration in Section 7.2. Specifically, we use the Generalized Synthetic Control Method (Xu, 2017) to estimate a factor model with three latent factors based on the series of log GSP per capita ($N = 50$, $T_0 = 89$). We then simulate outcomes using the distribution of estimated parameters and model selection into treatment as a function of the latent factors; see Appendix C for additional details. We also present results from three additional DGPs, each calibrated to estimates from the same data: (1) the factor model with quadruple the standard deviation of the noise term, (2) a unit and time fixed effects model, and (3) an autoregressive model with 3 lags.

We explore the role of augmentation using simple outcome estimators. For each DGP, we consider five estimators: (1) SCM alone, (2) ridge regression alone, (3) Ridge ASCM, (4) fixed effects alone, and (5) SCM augmented with fixed effects (i.e., de-meaned SCM), as shown in Equation (9). Figure 3 shows the Monte Carlo estimate of the absolute bias as a percentage of the absolute bias

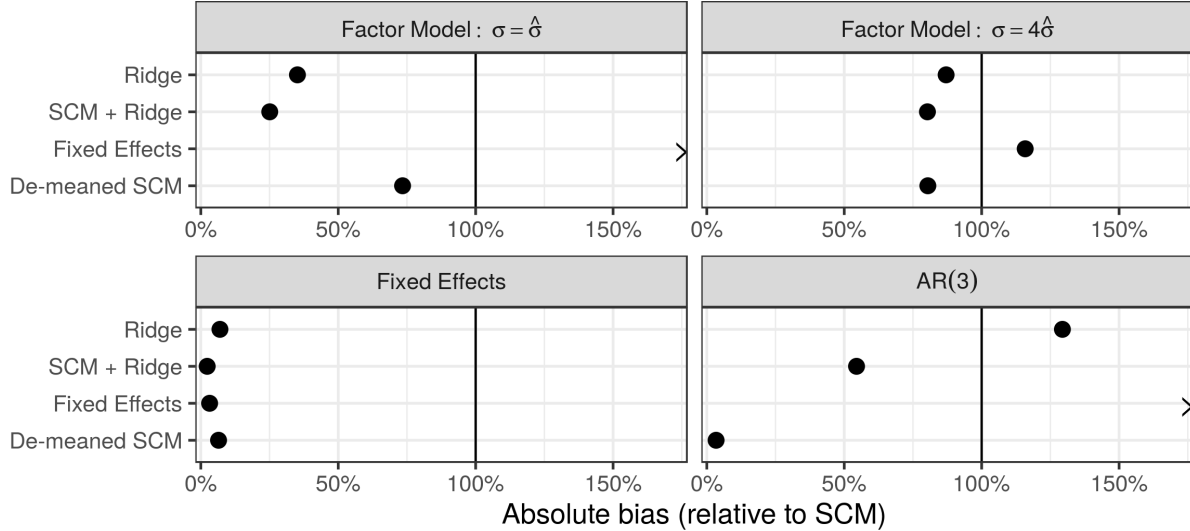


Figure 3: Overall absolute bias, normalized to SCM bias for (a) the factor model simulation, (b) the factor model simulation with quadruple the standard deviation, (c) the fixed effects simulation, and (d) the AR simulation. The SCM estimates reported here are *not* restricted to simulation draws with excellent pre-treatment fit; [Abadie et al. \(2015\)](#) advise against using SCM in such settings.

for SCM, with one panel for each simulation DGP; Appendix Figure F.10 shows the corresponding estimator root mean squared error (RMSE).

There are several takeaways. First, augmenting SCM with a ridge outcome regression reduces bias relative to SCM alone — *without* conditioning on excellent pre-treatment fit — in all four simulations. This underscores the importance of the recommendation in [Abadie et al. \(2010, 2015\)](#) to use SCM only in settings with excellent pre-treatment fit.⁶ Under the baseline factor model and the fixed effect model, the ridge augmentation greatly reduces bias, by more than 75% in the factor model simulation and over 90% in the fixed effects simulation. In the AR(3) model and in the factor model with greater noise, the gains to augmentation relative to SCM are more limited. Second, Ridge ASCM has slightly lower bias than ridge regression alone across all of the simulation settings. Third, when the fixed effects estimator is incorrectly specified, combining SCM with a fixed effects estimator has much lower bias than either method alone. And even when the fixed effects estimator is correctly specified, de-meaned SCM has similar bias to the (correctly specified) fixed effects approach. Finally, Appendix Figure F.10 shows that in all simulations ASCM has lower RMSE than SCM, as the large decrease in bias more than makes up for the slight increase in variance.

Complementing the worst-case analysis in Section 5, we now consider how the typical perfor-

⁶[Abadie et al. \(2010, 2015\)](#) also strongly recommend incorporating auxiliary covariates, weighted by their predictive power, into the procedure, noting that this is important for further reducing bias. For simplicity, the simulations do not include auxiliary covariates.

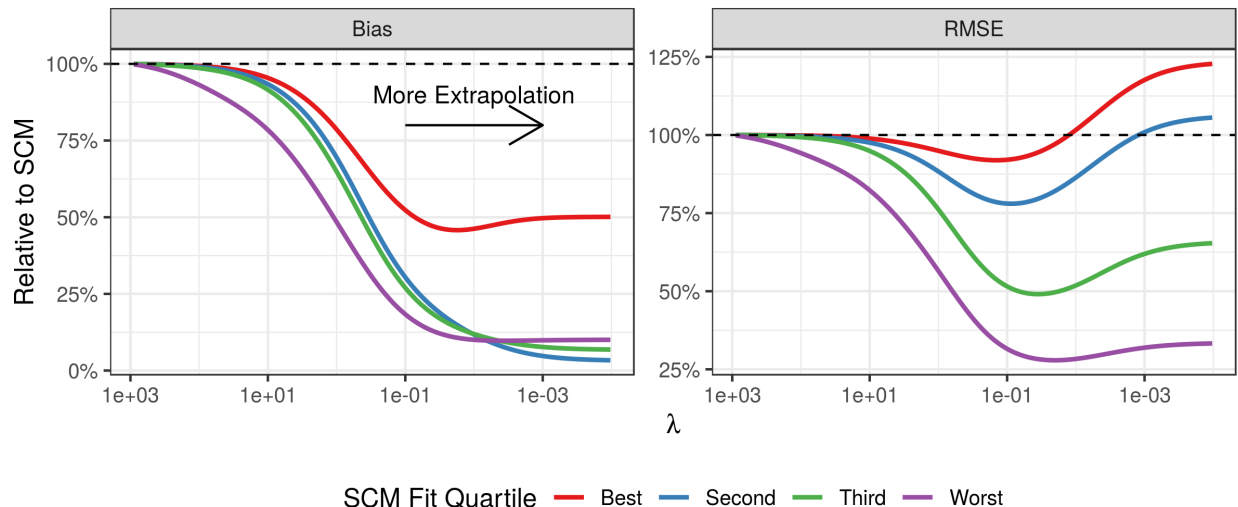


Figure 4: Bias and RMSE of Ridge ASCM, as a percentage of SCM bias and RMSE, versus λ under a linear factor model. Results are divided by the quartile of the SCM fit across all simulations.

mance of augmentation relates to the amount of extrapolation and the quality of the original SCM fit. Figure 4 shows the bias and RMSE as a function of λ for the primary factor model simulation, conditional on the quartile of SCM fit. Larger values of λ (and hence smaller adjustments) are to the left, with the left-most points in the plots representing SCM. First, as expected, Augmented SCM substantially reduces bias regardless of SCM pre-treatment fit. However, the gains are more modest when the SCM fit is in the best quartile: in this case the bias is non-monotonic in λ and there is some optimal choice of λ that minimizes the bias. Second, it is possible to under-regularize with ASCM, as evident in the RMSE achieving a minimum for an intermediate value of λ . Furthermore, when pre-treatment fit is good, augmentation with too-small λ leads to higher RMSE than SCM alone.

Next, we evaluate alternative outcome models for use in ASCM. For each DGP we consider SCM augmented with (1) LASSO, (2) a random forest, (3) `CausalImpact` (Brodersen et al., 2015), (4) matrix completion using `MCPanel` (Athey et al., 2017) and (5) fitting the factor model directly with `gsynth` (Xu, 2017). We compare ASCM to the pure outcome models as well as pure SCM. Figure 5 shows the absolute bias for these methods, again as a percentage of the absolute bias for SCM alone; Appendix Figure F.10 shows the RMSE. We broadly see the same results as with Ridge ASCM. In our simulations, augmenting SCM almost always reduces the bias relative to SCM (unconditional on good pre-treatment fit) with some models improving SCM more than others. Additionally, in nearly every case ASCM also has lower bias than outcome modeling alone. Appendix Figures F.11 and F.12 show the bias and RMSE when SCM fit is in the top quintile. As before, conditioned on good SCM pre-treatment fit the gains to augmentation with flexible outcome models are more

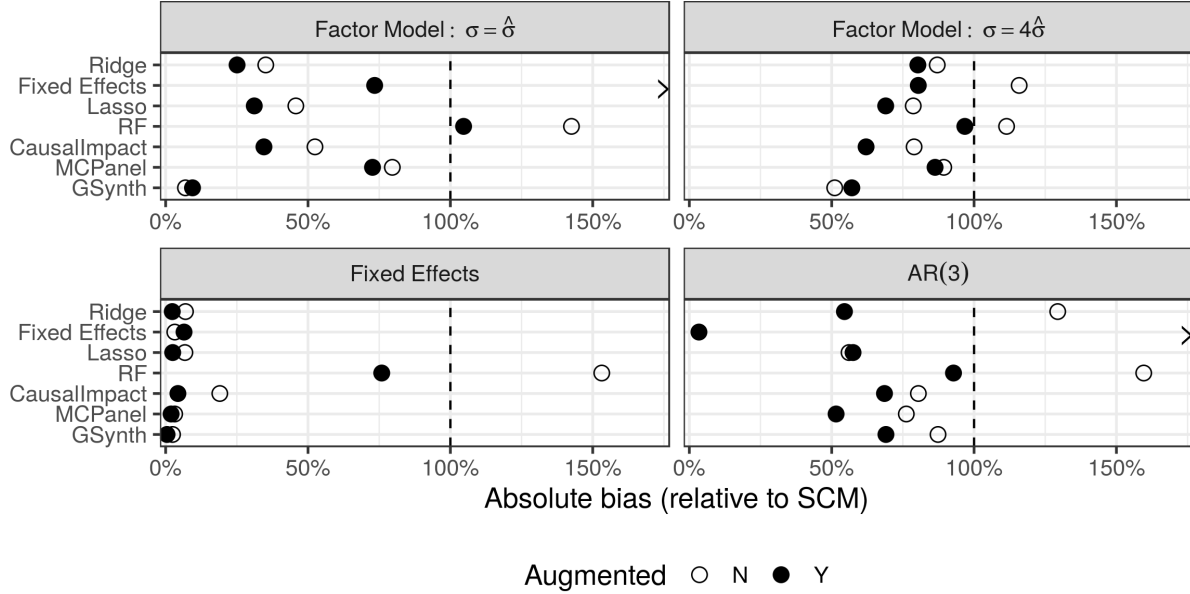


Figure 5: Absolute bias (as a percentage of SCM bias) for ridge, fixed effects, and several machine learning and panel data outcome models, and their augmented versions using the same data generating processes as Figure 3.

limited, except with the oracle `gsynth` estimator.

Overall we find that SCM augmented with a penalized regression model has consistently good performance across data generating processes. Due to this performance and the method’s relative simplicity, we therefore recommend augmenting SCM with penalized regression as a reasonable default in settings where SCM alone has poor pre-treatment fit. In particular, we suggest using ridge regression; among the other benefits, Ridge ASCM allows the practitioner to diagnose the level of extrapolation due to the outcome model.

7.2 Illustration: 2012 Kansas tax cuts

In 2010, Sam Brownback was elected governor of Kansas, having run on a platform emphasizing tax cuts and deficit reduction (see [Rickman and Wang, 2018](#), for further discussion and analysis). Upon taking office, he implemented a substantial personal income tax cut, both lowering rates and reducing credits and deductions. This is a valuable test of “supply side” models: Brownback argued that the tax cuts would increase business activity in Kansas, generating economic growth and additional tax revenues that would make up for the static revenue losses. Kansas’ subsequent economic performance has not been impressive relative to its neighbors; however, potentially confounding factors include a drought and declines in the locally important aerospace industry. Finding a credible control for Kansas is thus challenging, and SCM-type models offer a potential solution.

We estimate the effect of the tax cuts on log gross state product (GSP) per capita using the second quarter of 2012 — when Brownback signed the tax cut bill into law — as the intervention time. Results are consistent using outcomes scales other than the standard normalization of log GSP per capita (see Appendix F). We use four primary estimators: (1) SCM alone fit on the entire vector of lagged outcomes, (2) Ridge ASCM, (3) Ridge ASCM including auxiliary covariates in parallel to lagged outcomes and (4) Ridge ASCM on residualized outcomes, as in Section 6.⁷ These estimators rely on the ignorability assumption in Equation (2); substantively, this assumes that post-treatment shocks for Kansas will be the same as for other states in expectation. This also rules out unobserved confounders that affect both post-treatment shocks and the decision to enact the Brownback tax cut bill.

Figure 6, known as a “gap plot”, shows the difference between Kansas and its synthetic control before and after the passage of the tax cuts along with 95% point-wise confidence intervals computed via the conformal inference procedure from Chernozhukov et al. (2019); see Appendix A. Appendix F shows additional results, including estimates plotted on the raw outcome scale and results with alternative estimators. First, SCM alone achieves fairly poor pre-treatment fit; the synthetic control exceeds the treatment unit by two to four percent in 2004–2005, on the same scale as the estimated average post-treatment effect of a 3 percent decrease. This lack of pre-treatment fit should make us wary of the validity of the SCM effect estimates, and suggests that there may be gains to augmentation. Augmenting SCM with ridge regression indeed improves pre-treatment fit, especially in the mid 2000s. To better understand this, we can inspect the ridge regression coefficients for lagged outcomes (see Appendix Figure F.5), which put the most weight on the two most recent years. Adding auxiliary covariates and augmenting further improves both pre-treatment fit and balance on the covariates; see Figure 8a. Finally, balancing the auxiliary covariates via residualization also improves pre-treatment fit. Overall, the estimated impact is consistently negative for all four approaches, with weaker evidence that the effect persists to the end of the observation period.

To check against over-fitting, Figure 7 shows in-time placebo estimates for the Ridge ASCM estimator with covariates, with placebo treatment times in the second quarter of 2009, 2010, and 2011. We estimate placebo effects that are near zero with all three placebo treatment times. Appendix Figures F.6 and F.7 show the corresponding placebo estimates for SCM alone and Ridge ASCM without covariates.

Figure 8a shows the covariate balance for the four estimators. While SCM and Ridge ASCM achieve excellent fit for the pre-treatment average log GSP per capita, neither estimator achieves good balance on the other covariates, most notably the average employment level across the quarters of the pre-period. In contrast, including the auxiliary covariates into both the SCM and ridge

⁷The covariates we include are the pre-treatment averages of (1) log state and local revenue per capita, (2) log average weekly wages, (3) number of establishments per capita, (4) the employment level, and (5) log GSP per capita. For the augmented estimator on the lagged outcomes we select the hyperparameter λ^{ridge} as the largest λ within one standard error of the λ that minimizes the cross-validation placebo fit $CV(\lambda)$; see Section 5.3. Appendix Figure F.1 plots $CV(\lambda)$. When including the auxiliary covariates we use the minimal λ .

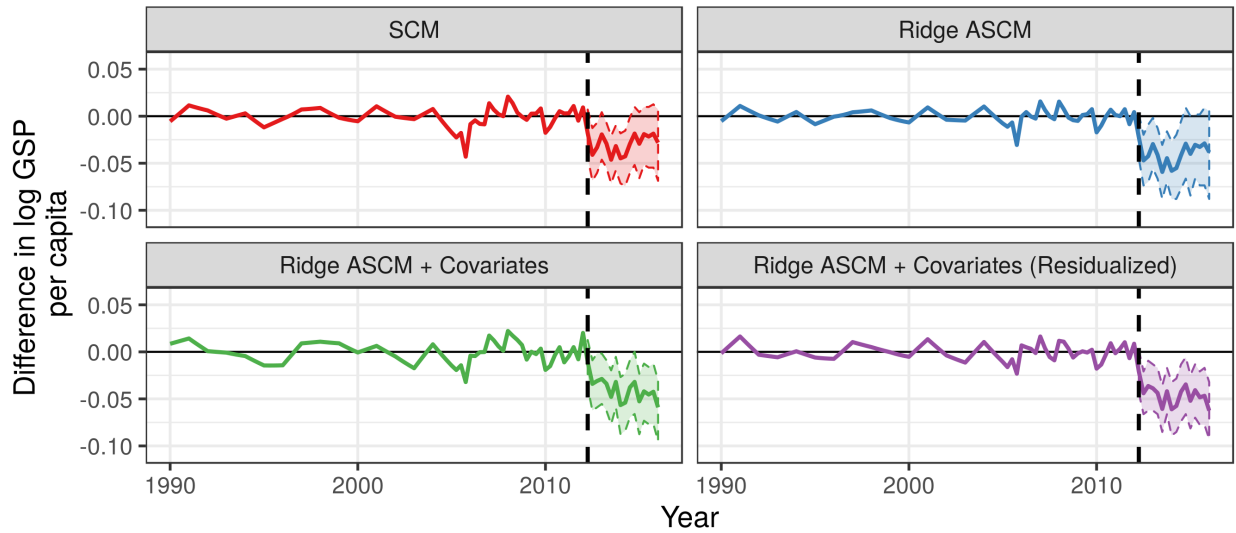


Figure 6: Point estimates along with point-wise 95% conformal confidence intervals for the effect of the tax cuts on log GSP per capita using SCM, Ridge ASCM, and Ridge ASCM with covariates.

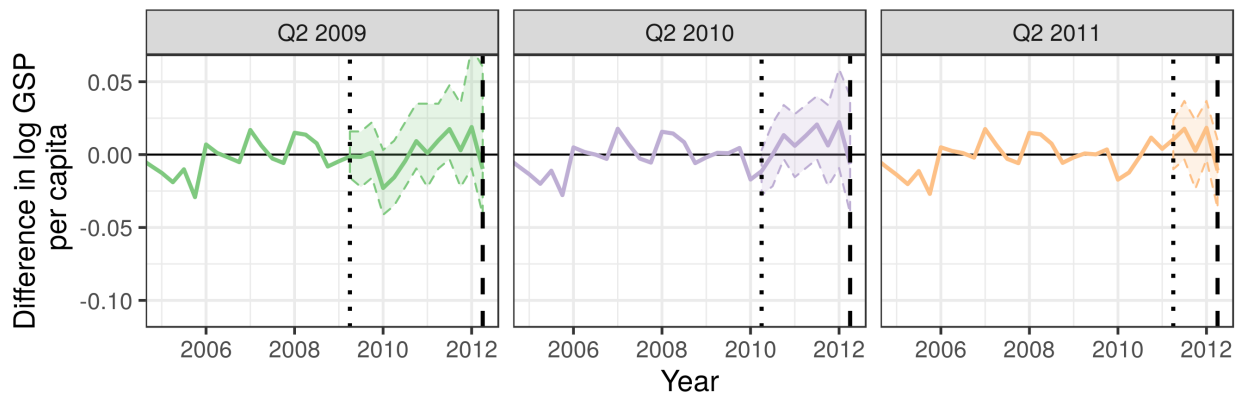


Figure 7: Placebo point estimates along with 95% conformal confidence intervals for Ridge ASCM with covariates with placebo treatment times in Q2 2009, 2010, and 2011. The time period begins in 2005 and ends in Q1 2012 to highlight placebo estimates.

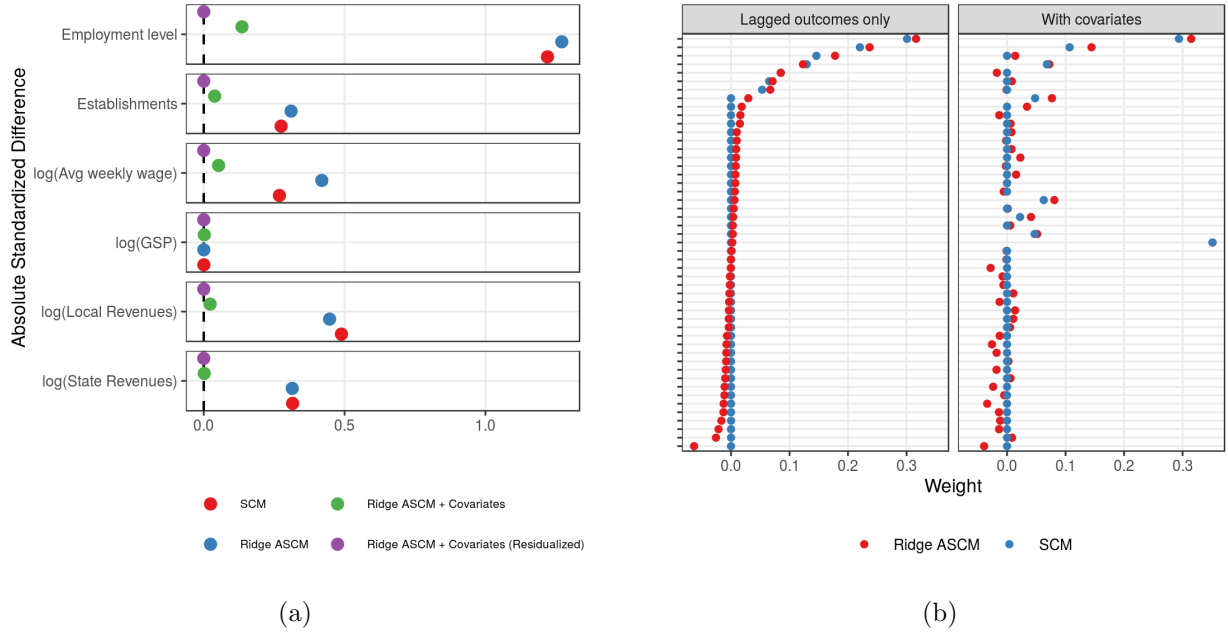


Figure 8: (a) Covariate balance for SCM, Ridge ASCM, and ASCM with covariates. Each covariate is standardized to have mean zero and standard deviation one; we plot the absolute difference between the treated unit’s covariate and the weighted control units’ covariates $|Z_{1k} - \sum_{W_i=0} \hat{\gamma} Z_{ik}|$. (b) Donor unit weights for SCM, Ridge ASCM, using lagged outcomes only or including auxiliary covariates.

optimization problems greatly improves the covariate balance, and — by design — residualizing on the auxiliary covariates perfectly balances them. Moreover, Ridge ASCM on residualized outcomes achieves very good pre-treatment fit on the lagged outcomes as shown in Figure 6.

Finally, Figure 8b shows the weights on donor units for SCM and Ridge ASCM as well as SCM and Ridge ASCM weights when including covariates jointly with the lagged outcomes. Appendix Figure F.9 shows the SCM and Ridge ASCM weights fit on the lagged outcomes after residualizing out the auxiliary covariates. Here we see the minimal extrapolation property of the ASCM weights. The SCM weights are zero for all but six donor states. The Ridge ASCM weights are similar but deviate slightly from the simplex, with only Louisiana receiving a meaningful negative weight. In addition, the Ridge ASCM weights retain some of the interpretability of the SCM weights. For the donor units with positive SCM weight, Ridge ASCM places close to the same weight. For the majority of those with zero SCM weight, Ridge ASCM also places a close to zero weight, with relatively few donor units with non-negligible negative weight. By contrast, Appendix Figure F.8 shows the weights from ridge regression alone: many of the weights are negative and the weights are far from sparse. Including auxiliary covariates changes the relative importance of different states by adding new information, but the minimal extrapolation property remains.

8 Discussion

SCM is a popular approach for estimating policy impacts at the jurisdiction level, such as the city or state. By design, however, the method is limited to settings where excellent pre-treatment fit is possible. For settings when this is infeasible, we introduce Augmented SCM, which controls pre-treatment fit while minimizing extrapolation. We show that this approach controls error under a linear factor model and propose several extensions, including to incorporate auxiliary covariates.

There are several directions for future work. The most immediate is to explore robust inferential methods for settings where pre-treatment SCM fit is imperfect. In Appendix A, we outline how to apply the conformal inference approach of [Chernozhukov et al. \(2019\)](#) to Augmented SCM, as well as a possible modification based on the jackknife+ approach of [Barber et al. \(2019\)](#). More work is needed, however, to understand the performance of this in more general settings. We could also extend this to allow for sensitivity analysis that directly parameterizes departures from, say, the linear factor model, as in recent approaches for sensitivity analysis for balancing weights ([Soriano et al., 2020](#)).

A second area for future inquiry is the application of the ASCM framework to settings with multiple treated units. For instance, there are different approaches in settings when all treated units are treated at the same time: some papers propose to fit SCM separately for each treated unit (e.g., [Abadie and L'Hour, 2018](#)), while others simply average the units together (e.g., [Kreif et al., 2016](#); [Robbins et al., 2017](#)). The situation is more complicated with staggered adoption, when units take up the treatment at different times (e.g., [Dube and Zipperer, 2015](#); [Donohue et al., 2017](#)). We explore this extension in [Ben-Michael et al. \(2019\)](#).

A third potential extension is to more complex data structures, such as applications with multiple outcomes series for the same units (e.g., measures of both earnings and total employment in minimum wage studies); hierarchical data structures with outcome information at both the individual and aggregate level (e.g., students within schools); or discrete or count outcomes.

References

- Abadie, A. (2019). Using synthetic controls: Feasibility, data requirements, and methodological aspects. *Journal of Economic Literature*.
- Abadie, A., A. Diamond, and J. Hainmueller (2010). Synthetic Control Methods for Comparative Case Studies: Estimating the Effect of California’s Tobacco Control Program. *Journal of the American Statistical Association* 105(490), 493–505.
- Abadie, A., A. Diamond, and J. Hainmueller (2015). Comparative Politics and the Synthetic Control Method. *American Journal of Political Science* 59(2), 495–510.
- Abadie, A. and J. Gardeazabal (2003). The Economic Costs of Conflict: A Case Study of the Basque Country. *The American Economic Review* 93(1), 113–132.
- Abadie, A. and G. W. Imbens (2011). Bias-corrected matching estimators for average treatment effects. *Journal of Business & Economic Statistics* 29(1), 1–11.
- Abadie, A. and J. L’Hour (2018). A penalized synthetic control estimator for disaggregated data.
- Amjad, M., D. Shah, and D. Shen (2018). Robust synthetic control. *The Journal of Machine Learning Research* 19(1), 802–852.
- Arkhangelsky, D., S. Athey, D. A. Hirshberg, G. W. Imbens, and S. Wager (2019). Synthetic difference in differences. *arXiv preprint arXiv:1812.09970*.
- Athey, S., M. Bayati, N. Doudchenko, G. Imbens, and K. Khosravi (2017). Matrix Completion Methods for Causal Panel Data Models. *arxiv 1710.10251*.
- Athey, S. and G. W. Imbens (2017). The state of applied econometrics: Causality and policy evaluation. *Journal of Economic Perspectives* 31(2), 3–32.
- Athey, S., G. W. Imbens, and S. Wager (2018). Approximate residual balancing: debiased inference of average treatment effects in high dimensions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 80(4), 597–623.
- Bai, J. (2009). Panel data models with interactive fixed effects. *Econometrica* 77(4), 1229–1279.
- Barber, R. F., E. J. Candès, A. Ramdas, and R. J. Tibshirani (2019). Predictive inference with the jackknife+. *arXiv preprint arXiv:1905.02928*.
- Ben-Michael, E., A. Feller, and J. Rothstein (2019). Synthetic controls and weighted event studies with staggered adoption. *arXiv preprint arXiv:1912.03290*.
- Bilinski, A. and L. Hatfield (2020). Goldilocks and the pre-intervention time series. Technical report.
- Botosaru, I. and B. Ferman (2019). On the role of covariates in the synthetic control method. *The Econometrics Journal* 22(2), 117–130.
- Breidt, F. J. and J. D. Opsomer (2017). Model-Assisted Survey Estimation with Modern Prediction Techniques. *Statistical Science* 32(2), 190–205.

- Brodersen, K. H., F. Gallusser, J. Koehler, N. Remy, and S. L. Scott (2015). Inferring Causal Impact using Bayesian Structural Time-Series Models. *The Annals of Applied Statistics* 9(1), 247–274.
- Cassel, C. M., C.-E. Sarndal, and J. H. Wretman (1976). Some results on generalized difference estimation and generalized regression estimation for finite populations. *Biometrika* 63(3), 615–620.
- Cattaneo, M. D., Y. Feng, and R. Titiunik (2019). Prediction intervals for synthetic control methods. *arXiv preprint arXiv:1912.07120*.
- Chattopadhyay, A., Christopher H. Hase, and J. R. Zubizarreta (2020). Balancing Versus Modeling Approaches to Weighting in Practice. *Statistics in Medicine in press*.
- Chernozhukov, V., K. Wuthrich, and Y. Zhu (2018). Inference on average treatment effects in aggregate panel data settings. *arXiv preprint arXiv:1812.10820*.
- Chernozhukov, V., K. Wüthrich, and Y. Zhu (2019). An Exact and Robust Conformal Inference Method for Counterfactual and Synthetic Controls. Technical report.
- Donohue, J. J., A. Aneja, and K. D. Weber (2017). Right-to-carry laws and violent crime: A comprehensive assessment using panel data and a state-level synthetic control analysis. Technical report, National Bureau of Economic Research.
- Doudchenko, N. and G. W. Imbens (2017). Difference-In-Differences and Synthetic Control Methods: A Synthesis. *arxiv 1610.07748*.
- Dube, A. and B. Zipperer (2015). Pooling multiple case studies using synthetic controls: An application to minimum wage policies.
- Ferman, B. (2019). On the Properties of the Synthetic Control Estimator with Many Periods and Many Controls.
- Ferman, B. and C. Pinto (2018). Synthetic controls with imperfect pre-treatment fit.
- Gobillon, L. and T. Magnac (2016). Regional policy evaluation: Interactive fixed effects and synthetic controls. *Review of Economics and Statistics* 98(3), 535–551.
- Hastie, T., J. Friedman, and R. Tibshirani (2009). *The elements of statistical learning*. Springer series in statistics New York.
- Hazlett, C. and Y. Xu (2018). Trajectory balancing: A general reweighting approach to causal inference with time-series cross-sectional data.
- Hirshberg, D. A., A. Maleki, and J. Zubizarreta (2019). Minimax linear estimation of the retargeted mean. *arXiv preprint arXiv:1901.10296*.
- Hirshberg, D. A. and S. Wager (2018). Augmented Minimax Linear Estimation.
- Hsiao, C., Q. Zhou, et al. (2018). Panel parametric, semi-parametric and nonparametric construction of counterfactuals-california tobacco control revisited. Technical report.

- Kellogg, M., M. Mogstad, G. Pouliot, and A. Torgovitsky (2020). Combining matching and synthetic controls to trade off biases from extrapolation and interpolation. Technical report, National Bureau of Economic Research.
- King, G. and L. Zeng (2006). The dangers of extreme counterfactuals. *Political Analysis* 14(2), 131–159.
- Kline, P. (2011). Oaxaca-Blinder as a reweighting estimator. In *American Economic Review*, Volume 101, pp. 532–537.
- Kreif, N., R. Grieve, D. Hangartner, A. J. Turner, S. Nikolova, and M. Sutton (2016). Examination of the synthetic control method for evaluating health policies with multiple treated units. *Health economics* 25(12), 1514–1528.
- Li, K. T. (2017). Statistical Inference for Average Treatment Effects Estimated by Synthetic Control Methods.
- Minard, S. and G. R. Waddell (2018). Dispersion-weighted synthetic controls.
- Neyman, J. (1990 [1923]). On the application of probability theory to agricultural experiments. essay on principles. section 9. *Statistical Science* 5(4), 465–472.
- Powell, D. (2018). Imperfect synthetic controls: Did the massachusetts health care reform save lives?
- Rickman, D. S. and H. Wang (2018). Two tales of two us states: Regional fiscal austerity and economic performance. *Regional Science and Urban Economics* 68, 46–55.
- Robbins, M., J. Saunders, and B. Kilmer (2017). A Framework for Synthetic Control Methods With High-Dimensional, Micro-Level Data: Evaluating a Neighborhood-Specific Crime Intervention. *Journal of the American Statistical Association* 112(517), 109–126.
- Robins, J. M., A. Rotnitzky, and L. P. Zhao (1994). Estimation of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association* 89(427), 846–866.
- Rubin, D. B. (1973). The use of matched sampling and regression adjustment to remove bias in observational studies. *Biometrics*, 185–203.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology* 66(5), 688.
- Rubin, D. B. (1980). Comment on “randomization analysis of experimental data: The fisher randomization test”. *Journal of the American Statistical Association* 75(371), 591–593.
- Samartsidis, P., S. R. Seaman, A. M. Presanis, M. Hickman, D. De Angelis, et al. (2019). Assessing the causal effect of binary interventions from observational panel data with few treated units. *Statistical Science* 34(3), 486–503.
- Soriano, D., E. Ben-Michael, P. Bickel, A. Feller, and S. Pimentel (2020). Sensitivity analysis for balancing weights. Technical report. working paper.

- Tan, Z. (2017). Regularized calibrated estimation of propensity scores with model misspecification and high-dimensional data.
- Toulis, P. and A. Shaikh (2018). Randomization tests in observational studies with time-varying adoption of treatment.
- Vovk, V., A. Gammerman, and G. Shafer (2005). *Algorithmic learning in a random world*. Springer.
- Wainwright, M. (2018). *High dimensional statistics: a non-asymptomatic viewpoint*.
- Wang, Y. and J. R. Zubizarreta (2018). Minimal Approximately Balancing Weights: Asymptotic Properties and Practical Considerations.
- Xu, Y. (2017). Generalized Synthetic Control Method: Causal Inference with Interactive Fixed Effects Models. *Political Analysis* 25, 57–76.
- Zhao, Q. (2018). Covariate Balancing Propensity Score by Tailored Loss Functions. *Annals of Statistics*, forthcoming.
- Zhao, Q. and D. Percival (2017). Entropy balancing is doubly robust. *Journal of Causal Inference* 5(1).
- Zubizarreta, J. R. (2015). Stable Weights that Balance Covariates for Estimation With Incomplete Outcome Data. *Journal of the American Statistical Association* 110(511), 910–922.

Supplementary Materials for “The Augmented Synthetic Control Method”

A Inference

There is a large and growing literature on inference for the synthetic control method and variants, going beyond the original proposal in [Abadie and Gardeazabal \(2003\)](#) and [Abadie et al. \(2010\)](#); see, for example, [Li \(2017\)](#), [Toulis and Shaikh \(2018\)](#), [Cattaneo et al. \(2019\)](#), [Toulis and Shaikh \(2018\)](#), and [Chernozhukov et al. \(2018\)](#). Here, we consider the conformal inference approach of [Chernozhukov et al. \(2019\)](#), which is tailored to this setting, as well as an adaptation of the jackknife+ approach of [Barber et al. \(2019\)](#).

We now briefly describe the conformal inference approach of [Chernozhukov et al. \(2019\)](#):

1. For a given sharp null hypothesis, $H_0 : \tau = \tau_0$:
 - (a) Enforce the null hypothesis by creating an adjusted post-treatment outcome for the treated unit $\tilde{Y}_{1T} = Y_{1T} - \tau_0$.
 - (b) Augment the original data set to include the post-treatment time period T , with the adjusted outcome \tilde{Y}_{1T} ; use the estimator (13) to obtain adjusted weights $\hat{\gamma}(\tau_0)$.
 - (c) Compute a p -value by assessing whether the adjusted residual $Y_{1T} - \tau_0 - \sum_{W_i=0} \hat{\gamma}_i(\tau_0) Y_{iT}$ “conforms” with the pre-treatment residuals:⁸

$$\hat{p}(\tau_0) = \frac{1}{T} \sum_{t=1}^{T_0} \mathbb{1} \left\{ \left| Y_{1T} - \tau_0 - \sum_{W_i=0} \hat{\gamma}_i(\tau_0) Y_{iT} \right| \leq \left| Y_{1t} - \sum_{W_i=0} \hat{\gamma}_i(\tau_0) Y_{it} \right| \right\} + \frac{1}{T}. \quad (\text{A.1})$$

2. Compute a level α confidence interval for τ by inverting the hypothesis test and constructing the set

$$\hat{C}_\tau^{\text{conf}}(\alpha) = \{\tau_0 \mid \hat{p}_{\tau_0} \geq \alpha\}.$$

Since the counterfactual outcome $Y_{1T}(0)$ is random, this is equivalent to constructing a conformal prediction set ([Vovk et al., 2005](#)) for $Y_{1T}(0)$ by using the quantiles of pre-treatment residuals:

$$\hat{C}_Y^{\text{conf}} = \left\{ y \in \mathbb{R} \mid \left| y - \sum_{W_i=0} \hat{\gamma}_i(Y_{1T} - y) Y_{iT} \right| \leq q_{T,\alpha}^+ \left(\left| Y_{1t} - \sum_{W_i=0} \hat{\gamma}_i(Y_{1T} - y) Y_{it} \right| \right) \right\}, \quad (\text{A.2})$$

where $q_{T,\alpha}^+(x_t)$ is the $[(1 - \alpha)T]^{\text{th}}$ order statistic of \mathbf{X}_1, \dots, x_T . Note that $\tau \in \hat{C}_\tau^{\text{conf}} \Leftrightarrow Y_{1T}(0) \in \hat{C}_Y^{\text{conf}}$. If $\tau \in \hat{C}_\tau^{\text{conf}}$, then the adjusted residual is less than or equal to the $[(1 - \alpha)T]^{\text{th}}$ smallest pre-treatment residual and so $Y_{1T}(0) \in \hat{C}_Y^{\text{conf}}$. Conversely if $\tau \notin \hat{C}_\tau^{\text{conf}}$, then the adjusted residual must be larger than the $[(1 - \alpha)T]^{\text{th}}$ smallest pre-treatment residual and so $Y_{1T}(0) \notin \hat{C}_Y^{\text{conf}}$.

[Chernozhukov et al. \(2019\)](#) provide several conditions for which approximate or exact finite-sample validity of the p -values (and hence coverage of the prediction interval \hat{C}_Y^{conf}) can be achieved;

⁸There are several choices, such as the test statistic and the form of permutation across time periods, that reduce to Equation (A.1) with a single post treatment time period. See [Chernozhukov et al. \(2019\)](#) for further details.

our setup in Section 2.1 follows theirs. First, they show that approximate validity under an additive noise model with either i.i.d. or stationary noise depends on the estimation error. Applying Proposition 1 and Theorem 1, we can characterize the finite-sample coverage probability of $\widehat{C}_Y^{\text{conf}}$ and see that the true coverage will be close to the nominal coverage level α if the pre-treatment fit is good and the approximation error is small.

Second, they show exact validity if the residuals are exchangeable. Intuitively, pre-treatment residuals must be a good proxy for post-treatment residuals; under the linear factor model this can hold when the factor values do not differ much across time periods (i.e., for a two way fixed effects model), but can be violated if the factor values differ widely in the pre- and post-intervention periods. Additionally, if the ignorability assumption does not hold and treatment adoption is correlated with the shocks ε_{it} , then pre-treatment residuals will be poor proxies for post-treatment residuals, leading to undercoverage.

One drawback of the “full” conformal approach is that it is computationally intensive: to construct the prediction interval $\widehat{C}_Y^{\text{conf}}$, we must refit the weights $\hat{\gamma}_i^{\text{aug}}(Y_{1T} - y)$ for a grid of possible values of the true counterfactual outcome y . A recent alternative conformal approach is the jackknife+ (Barber et al., 2019). This procedure uses the leave one out residuals $Y_{1t} - \hat{Y}_{1t}^{(-t)}$ and the estimate of the post-treatment period $\hat{Y}_{1T}^{(-t)}$ after dropping period t . The prediction interval for $Y_{1T}(0)$ is

$$\widehat{C}_Y^{\text{jackknife+}} = \left[q_{T,\alpha/2}^- \left(\hat{Y}_{1T}^{(-t)} - \left| Y_{1t} - \hat{Y}_{1t}^{(-t)} \right| \right), q_{T,\alpha/2}^+ \left(\hat{Y}_{1T}^{(-t)} + \left| Y_{1t} - \hat{Y}_{1t}^{(-t)} \right| \right) \right], \quad (\text{A.3})$$

where $q_{T,\alpha}^-(x_t)$ is the $[\alpha T]^{\text{th}}$ order statistic of \mathbf{X}_1, \dots, x_T . As with the full conformal method, $\widehat{C}_Y^{\text{jackknife+}}$ will have exact finite sample coverage when the time periods or residuals are exchangeable. However, the jackknife+ procedure only requires that we re-fit the estimator for each of the T_0 pre-treatment time periods. While the full conformal method enforces a sharp null hypothesis when estimating the weights, the jackknife+ uses the leave-one-out residuals as a proxy for the distribution of $Y_{1T}(0)$, incorporating variability in the estimator by including the leave-one-out estimate of the post-treatment outcome $\hat{Y}_{1T}^{(-t)}$. We anticipate that it is possible to extend the results of Chernozhukov et al. (2019) to show approximate validity when residuals are not exchangeable; we leave this to future work.

Simulation study. We assess the finite sample coverage of the conformal prediction intervals for $Y_{1T}(0)$ using both the full conformal method (A.1) and the jackknife+ procedure (A.3). For the four simulation settings we compute 95% prediction intervals for the first post-treatment counterfactual outcome $Y_{1T_0+1}(0)$ using both conformal methods and the both the SCM and ridge ASCM estimators. Table A.1 shows the results. We see that the intervals for SCM sometimes undercover, due to finite sample bias from poor treatment fit. In contrast, the intervals for ridge ASCM have greater than or close to nominal coverage for Y_{1T_0+1} .

Model	Estimation method	Inference method	Coverage
AR(3)	SCM	Full conformal	0.934
		Jackknife+	0.950
	SCM + Ridge	Full conformal	0.932
		Jackknife+	0.947
Factor Model	SCM	Full conformal	0.926
		Jackknife+	0.954
	SCM + Ridge	Full conformal	0.950
		Jackknife+	0.966
Factor Model (More Noise)	SCM	Full conformal	0.930
		Jackknife+	0.956
	SCM + Ridge	Full conformal	0.936
		Jackknife+	0.957
Fixed Effects	SCM	Full conformal	0.889
		Jackknife+	0.957
	SCM + Ridge	Full conformal	0.939
		Jackknife+	0.956

Table A.1: Coverage for the full conformal (A.1) and jackknife+ (A.3) prediction intervals.

B Additional results

B.1 Specialization of Ridge ASCM results to SCM

This appendix section specializes select results from the main text for Ridge ASCM for the special case of SCM, with $\lambda \rightarrow \infty$.

First we specialize Proposition 1 to SCM weights by taking $\lambda \rightarrow \infty$.

Corollary A.1. Under the linear model (18) with independent sub-Gaussian noise with scale parameter σ , for any $\delta > 0$, for weights $\gamma \in \Delta^{N_0}$ independent of the post-treatment outcomes (Y_{1T}, \dots, Y_{NT}) and for any $\delta > 0$,

$$Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i Y_{iT} \leq \underbrace{\|\beta\|_2 \left\| \mathbf{X}_1 - \sum_{W_i=0} \hat{\gamma}_i \mathbf{X}_i \right\|_2}_{\text{imbalance in } \mathbf{X}} + \underbrace{\delta \sigma (1 + \|\hat{\gamma}\|_2)}_{\text{post-treatment noise}}, \quad (\text{A.4})$$

with probability at least $1 - 2e^{-\frac{\delta^2}{2}}$.

We can similarly specialize Theorem 1.

Corollary A.2. Under the linear factor model (21) with independent sub-Gaussian noise with scale parameter σ , for weights $\gamma \in \Delta^{N_0}$ independent of the post-treatment outcomes (Y_{1T}, \dots, Y_{NT}) and for any $\delta > 0$,

$$Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i Y_{iT} \leq \underbrace{\frac{JM^2}{\sqrt{T_0}} \left\| \mathbf{X}_1 - \sum_{W_i=0} \hat{\gamma}_i \mathbf{X}_i \right\|_2}_{\text{imbalance in } \mathbf{X}} + \underbrace{\frac{2JM^2\sigma}{\sqrt{T_0}} (\sqrt{\log 2N_0} + \delta)}_{\text{approximation error}} + \underbrace{\delta \sigma (1 + \|\hat{\gamma}\|_2)}_{\text{post-treatment noise}}, \quad (\text{A.5})$$

with probability at least $1 - 6e^{-\frac{\delta^2}{2}}$.

B.2 Error under a partially linear model with Lipschitz deviations from linearity

We now bound the estimation error for SCM and Ridge ASCM when the outcome is only partially linear, with Lipschitz deviations from linearity. Specifically, assume that the control potential outcome in time period T satisfies

$$Y_{iT}(0) = \beta \cdot \mathbf{X}_i + f(\mathbf{X}_i) + \varepsilon_{iT}, \quad (\text{A.6})$$

where $f : \mathbb{R}^{T_0} \rightarrow \mathbb{R}$ is L -Lipschitz and the noise terms ε_{it} are independent sub-Gaussian random variables with scale parameter σ , and are ignorable $\mathbb{E}_{\varepsilon_T}[W_i \varepsilon_{iT}] = \mathbb{E}_{\varepsilon_T}[(1 - W_i) \varepsilon_{iT}] = \mathbb{E}_{\varepsilon_T}[\varepsilon_{iT}] = 0$, as above.

Under this model, the L -Lipschitz function $f(\cdot)$ will induce an approximation error from deviating away from the nearest neighbor match.

Theorem A.1. Let $C = \max_{W_i=0} \|\mathbf{X}_i\|_2$. For any $\delta > 0$, the estimation error for the ridge ASCM weights $\hat{\gamma}^{\text{aug}}$ (13) with hyperparameter $\lambda^{\text{ridge}} = N_0 \lambda$ is

$$\begin{aligned}
\left| Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i^{\text{aug}} Y_{1T} \right| &\leq \underbrace{\|\beta\|_2 \left\| \text{diag} \left(\frac{\lambda}{d_j^2 + \lambda} \right) (\tilde{\mathbf{X}}_1 - \tilde{\mathbf{X}}_0' \hat{\gamma}^{\text{scm}}) \right\|_2}_{\text{imbalance in } X} + \\
&\quad \underbrace{CL \left\| \text{diag} \left(\frac{d_j}{d_j^2 + \lambda} \right) (\tilde{\mathbf{X}}_1 - \tilde{\mathbf{X}}_0' \hat{\gamma}^{\text{scm}}) \right\|_2}_{\text{excess approximation error}} + \\
&\quad \underbrace{L \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} \|\mathbf{X}_1 - \mathbf{X}_i\|_2}_{\text{SCM approximation error}} + \underbrace{\delta\sigma (1 + \|\hat{\gamma}^{\text{aug}}\|_2)}_{\text{post-treatment noise}}
\end{aligned} \tag{A.7}$$

with probability at least $1 - 2e^{-\frac{\delta^2}{2}}$.

We can again specialize this to the SCM weights alone by taking $\lambda \rightarrow \infty$.

Corollary A.3. For any $\delta > 0$, the estimation error for weights on the simplex $\hat{\gamma} \in \Delta^{N_0}$ independent of the post-treatment outcomes (Y_{1T}, \dots, Y_{NT}) is

$$\begin{aligned}
Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i Y_i &\leq \underbrace{\|\beta\|_2 \left\| \mathbf{X}_1 - \sum_{W_i=0} \hat{\gamma}_i \mathbf{X}_i \right\|_2}_{\text{imbalance in } X} + \underbrace{L \sum_{W_i=0} \hat{\gamma}_i \|\mathbf{X}_1 - \mathbf{X}_i\|_2}_{\text{approximation error}} + \underbrace{\delta\sigma(1 + \|\hat{\gamma}\|_2)}_{\text{post-treatment noise}}
\end{aligned} \tag{A.8}$$

with probability at least $1 - 2e^{-\frac{\delta^2}{2}}$.

Inspecting Corollary A.3, we see that in order to control the estimation error, the weights must ensure good pre-treatment fit, while only weighting control units that are near to the treated unit, with the ratio $L/\|\beta\|_2$ controlling the relative importance of both terms. [Abadie and L'Hour \(2018\)](#) propose finding weights by solving the penalized SCM problem,

$$\min_{\gamma \in \Delta^{N_0}} \left\| \mathbf{X}_1 - \sum_{W_i=0} \hat{\gamma}_i \mathbf{X}_i \right\|_2^2 + \lambda \sum_{W_i=0} \hat{\gamma}_i \|\mathbf{X}_1 - \mathbf{X}_i\|_2^2. \tag{A.9}$$

Comparing this to Corollary A.3, we see that under the partially linear model (A.6) where $f(\cdot)$ is L -Lipshitz, finding weights that limit interpolation error by controlling both the overall imbalance in the lagged outcomes as well as the weighted sum of the distances is sufficient to control the error. In the above optimization problem, the hyperparameter λ takes the role of $L/\|\beta\|_2$.

B.3 Error under a linear factor model with covariates

We can quantify the behavior of the two-step procedure from Lemma 4 in controlling the error under a more general form of the linear factor model (21) with covariates (see Abadie et al., 2010; Botosaru and Ferman, 2019, for additional discussion). We can also consider the error under a linear model with auxiliary covariates, as a direct consequence of Lemma 4.

For each time period t , the covariates in the linear factor model enter through a time-varying function $f_t : \mathbb{R}^K \rightarrow \mathbb{R}$ in the model for the outcomes at Y_{it} :

$$Y_{it}(0) = \sum_{j=1}^J \phi_{ij} \mu_{jt} + f_t(\mathbf{Z}_i) + \varepsilon_{it}. \quad (\text{A.10})$$

For this model we again assume that ε_{iT} are independent, mean-zero sub-Gaussian random variables with scale parameter σ , and are uncorrelated with treatment assignment.

To characterize how well the covariates approximate the true function $f(\mathbf{Z}_i)$, we will consider the best linear approximation in our data, and define the residual for unit i and time t as $e_{it} = f_t(\mathbf{Z}_i) - \mathbf{Z}'_i(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'f_t(\mathbf{Z})$, where $\mathbf{Z} \in \mathbb{R}^{N \times K}$ is the matrix of all auxiliary covariates for all units. For each time period we will characterize the additional approximation error incurred by only balancing the covariates linearly with the *residual sum of squares* $RSS_t = \sum_{i=1}^n e_{it}^2$. For ease of exposition, we assume that the control covariates are standardized and rotated, which can always be true after pre-processing, and present results for the simpler case in which we fit SCM on the residualized pre-treatment outcomes rather than ridge ASCM (i.e., we let $\lambda^{\text{ridge}} \rightarrow \infty$); the more general version follows immediately by applying Theorem 1.

Theorem A.2. Under the linear factor model with covariates (A.12) with $\frac{1}{N_0}\mathbf{Z}'_0\mathbf{Z}_0 = \mathbf{I}_K$ and sub-Gaussian noise, for any $\delta > 0$, $\hat{\gamma}^{\text{cov}}$ in Equation (29) with $\lambda^{\text{ridge}} \rightarrow \infty$ satisfies the bound

$$\begin{aligned} \left| Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}^{\text{cov}} Y_{iT} \right| &\leq \frac{JM^2}{\sqrt{T_0}} \left(\underbrace{\|\check{\mathbf{X}}_1 - \check{\mathbf{X}}'_0 \hat{\gamma}\|_2}_{\text{imbalance in } \check{\mathbf{X}}} + 4\sigma \underbrace{\sqrt{\frac{K}{N_0}} \|\mathbf{Z}_1 - \mathbf{Z}'_0 \hat{\gamma}\|_2}_{\text{excess approximation error}} \right) + \\ &\quad \underbrace{\frac{2JM^2\sigma}{\sqrt{T_0}} \left(\sqrt{\log N_0} + \frac{\delta}{2} \right)}_{\text{SCM approximation error}} + \underbrace{(JM^2 + 1)e_{1\max} + (JM^2 + 1)\sqrt{RSS_{\max}} \|\hat{\gamma}^{\text{cov}}\|_2}_{\text{covariate approximation error}} \\ &\quad + \underbrace{\delta\sigma(1 + \|\hat{\gamma}^{\text{cov}}\|_2)}_{\text{post-treatment noise}} \end{aligned} \quad (\text{A.11})$$

with probability at least $1 - 6e^{-\frac{\delta^2}{2}} - 2e^{-\frac{KN_0(2-\sqrt{\log 5})^2}{2}}$, where $e_{1\max} = \max_t |e_{1t}|$ is the maximal residual for the treated unit and $RSS_{\max} = \max_t RSS_t$ is the maximal residual sum of squares

We can also consider the special case of Theorem A.2 when $f_t(\mathbf{Z}_i) = \sum_{k=1}^K B_{tk} Z_{ik}$ is a linear function of the covariates, and so

$$Y_{it}(0) = \sum_{j=1}^J \phi_{ij} \mu_{jt} + \sum_{k=1}^K B_{tk} Z_{ik} + \varepsilon_{it} = \boldsymbol{\phi}'_i \boldsymbol{\mu}_T + \mathbf{B}'_t \mathbf{Z}_i + \varepsilon_{it}. \quad (\text{A.12})$$

In this case the residuals $e_{it} = 0 \quad \forall i, t$.

Corollary A.4. Under the linear factor model with covariates (A.12) with $\frac{1}{N_0} \mathbf{Z}'_0 \mathbf{Z}_0 = \mathbf{I}_K$ and sub-Gaussian noise, for any $\delta > 0$, $\hat{\gamma}^{\text{cov}}$ in Equation (29) with $\lambda^{\text{ridge}} \rightarrow \infty$ satisfies the bound

$$\begin{aligned} \left| Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}^{\text{cov}} Y_{iT} \right| &\leq \frac{JM^2}{\sqrt{T_0}} \left(\underbrace{\|\check{\mathbf{X}}_1 - \check{\mathbf{X}}'_0 \hat{\gamma}\|_2}_{\text{imbalance in } \check{\mathbf{X}}} + \underbrace{4\sigma \sqrt{\frac{K}{N_0}} \|\mathbf{Z}_1 - \mathbf{Z}'_0 \hat{\gamma}\|_2}_{\text{excess approximation error}} \right) + \\ &\quad \underbrace{\frac{2JM^2\sigma}{\sqrt{T_0}} \left(\sqrt{\log N_0} + \frac{\delta}{2} \right)}_{\text{SCM approximation error}} + \underbrace{\delta\sigma(1 + \|\hat{\gamma}^{\text{cov}}\|_2)}_{\text{post-treatment noise}} \end{aligned} \tag{A.13}$$

with probability at least $1 - 6e^{-\frac{\delta^2}{2}} - 2e^{-\frac{KN_0(2-\sqrt{\log 5})^2}{2}}$.

Building on Lemma 4, Theorem A.2 and Corollary A.4 show that due to the additive, separable structure of the auxiliary covariates in Equation (A.12), controlling the pre-treatment fit in the *residualized* lagged outcomes $\check{\mathbf{X}}$ partially controls the error. This justifies directly targeting fit in the residualized lagged outcomes $\check{\mathbf{X}}$ rather than targeting raw lagged outcomes \mathbf{X} . Moreover, the excess approximation error will be small since since the number of covariates K is small relative to N_0 and the auxiliary covariates are measured without noise. As in Section 4.2, we can achieve better balance when we apply ridge ASCM to $\check{\mathbf{X}}$ than when we apply SCM alone. Because $\check{\mathbf{X}}$ are orthogonal to \mathbf{Z} by construction, this comes at no cost in terms of imbalance in \mathbf{Z} . However, the fundamental challenge of over-fitting to noise still remains, and, as in the case without auxiliary covariates, selecting the tuning parameter remains important. We again propose to follow the cross validation approach in Section 5.3, here using the residualized lagged outcomes $\check{\mathbf{X}}$ rather than the raw lagged outcomes \mathbf{X} .

C Simulation data generating process

We now describe the simulations in detail. We use the Generalized Synthetic Control Method (Xu, 2017) to fit the following linear factor model to the observed series of log GSP per capita ($N = 50, T_0 = 89, T = 105$), setting $J = 3$:

$$Y_{it} = \alpha_i + \nu_t + \sum_{j=1}^J \phi_{ij} \mu_{jt} + \varepsilon_{it}. \quad (\text{A.14})$$

We then use these estimates as the basis for simulating data. Appendix Figure F.13 shows the estimated factors $\hat{\boldsymbol{\mu}}$. We use the estimated time fixed effects $\hat{\boldsymbol{\nu}}$ and factors $\hat{\boldsymbol{\mu}}$ and then simulate data using Equation (A.14), drawing:

$$\begin{aligned} \alpha_i &\sim N(\hat{\alpha}, \hat{\sigma}_\alpha) \\ \phi &\sim N(0, \hat{\boldsymbol{\Sigma}}_\phi) \\ \varepsilon_{it} &\sim N(0, \hat{\sigma}_\varepsilon), \end{aligned}$$

where $\hat{\alpha}$ and $\hat{\sigma}_\alpha$ are the estimated mean and standard deviation of the unit-fixed effects, $\hat{\boldsymbol{\Sigma}}_\phi$ is the sample covariance of the estimated factor loadings, and $\hat{\sigma}_\varepsilon$ is the estimated residual standard deviation. We also simulate outcomes with quadruple the standard deviation, $\text{sd}(\varepsilon_{it}) = 4\hat{\sigma}_\varepsilon$. We assume a sharp null of zero treatment effect in all DGPs and estimate the ATT at the final time point.

To model selection, we compute the (marginal) propensity scores as

$$\text{logit}^{-1}\{\pi_i\} = \text{logit}^{-1}\{\mathbb{P}(T = 1 \mid \alpha_i, \boldsymbol{\phi}_i)\} = \theta \left(\alpha_i + \sum_j \phi_{ij} \right),$$

where we set $\theta = 1/2$ and re-scale the factors and fixed effects to have unit variance. Finally, we restrict each simulation to have a single treated unit and therefore normalize the selection probabilities as $\frac{\pi_i}{\sum_j \pi_j}$.

We also consider an alternative data generating process that specializes the linear factor model to only include unit- and time-fixed effects:

$$Y_{it}(0) = \alpha_i + \nu_t + \varepsilon_{it}.$$

We calibrate this data generating process by fitting the fixed effects with `gsynth` and drawing new unit-fixed effects from $\alpha_i \sim N(\hat{\alpha}, \hat{\sigma}_\alpha)$. We then model selection proportional to the fixed effect as above with $\theta = \frac{3}{2}$. Second, we generate data from an AR(3) model:

$$Y_{it}(0) = \beta_0 + \sum_{j=1}^3 \beta_j Y_{i(t-j)} + \varepsilon_{it},$$

where we fit $\beta_0, \boldsymbol{\beta}$ to the observed series of log GSP per capita. We model selection as proportional to the last 3 outcomes $\text{logit}^{-1}\pi_i = \theta \left(\sum_{j=1}^3 Y_{i(T_0-j+1)} \right)$ and set $\theta = \frac{5}{2}$. For this simulation we

estimate the ATT at time $T_0 + 1$.

D Proofs

D.1 Proofs for Section 4

Lemma A.1. With $\hat{\eta}_0^{\text{ridge}}$ and $\hat{\boldsymbol{\eta}}^{\text{ridge}}$, the solutions to (10), the ridge estimate can be written as a weighting estimator:

$$\hat{Y}_{1T}^{\text{ridge}}(0) = \hat{\eta}_0^{\text{ridge}} + \hat{\boldsymbol{\eta}}^{\text{ridge}'} \mathbf{X}_1 = \sum_{W_i=0} \hat{\gamma}_i^{\text{ridge}} Y_{iT}, \quad (\text{A.15})$$

where

$$\hat{\gamma}_i^{\text{ridge}} = \frac{1}{N_0} + (\mathbf{X}_1 - \bar{X}_0)' (\mathbf{X}'_0 \mathbf{X}_0 + \lambda^{\text{ridge}} \mathbf{I}_{T_0})^{-1} \mathbf{X}_i. \quad (\text{A.16})$$

Moreover, the ridge weights $\hat{\gamma}^{\text{ridge}}$ are the solution to

$$\boldsymbol{\gamma} \mid \sum_i \gamma_i = 1 \quad \min \frac{1}{2\lambda^{\text{ridge}}} \|\mathbf{X}_1 - \mathbf{X}'_0 \boldsymbol{\gamma}\|_2^2 + \frac{1}{2} \left\| \boldsymbol{\gamma} - \frac{1}{N_0} \mathbf{1} \right\|_2^2. \quad (\text{A.17})$$

Proof of Lemmas 1 and A.1. Recall that the lagged outcomes are centered by the control averages. Notice that

$$\begin{aligned} \hat{Y}_{1T}^{\text{aug}}(0) &= \hat{m}(\mathbf{X}_1) + \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} (Y_{iT} - \hat{m}(\mathbf{X}_i)) \\ &= \hat{\eta}_0 + \hat{\boldsymbol{\eta}}' \mathbf{X}_1 + \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} (Y_{iT} - \hat{\eta}_0 - \mathbf{X}'_i \hat{\boldsymbol{\eta}}) \\ &= \sum_{W_i=0} (\hat{\gamma}_i^{\text{scm}} + (\mathbf{X}_1 - \mathbf{X}'_0 \hat{\boldsymbol{\gamma}}^{\text{scm}}) (\mathbf{X}'_0 \mathbf{X}_0 + \lambda \mathbf{I}_{T_0})^{-1} \mathbf{X}_i) Y_{iT} \\ &= \sum_{W_i=0} \hat{\gamma}_i^{\text{aug}} Y_{iT} \end{aligned} \quad (\text{A.18})$$

The expression for $\hat{Y}_{1T}^{\text{ridge}}(0)$ follows.

We now prove that $\hat{\boldsymbol{\gamma}}^{\text{ridge}}$ and $\hat{\boldsymbol{\gamma}}^{\text{scm}}$ solve the weighting optimization problems (A.17) and (14). First, the Lagrangian dual to (A.17) is

$$\min_{\alpha, \boldsymbol{\beta}} \frac{1}{2} \sum_{W_i=0} \left(\alpha + \boldsymbol{\beta}' \mathbf{X}_i + \frac{1}{N_0} \right)^2 - (\alpha + \boldsymbol{\beta}' \mathbf{X}_1) + \frac{\lambda}{2} \|\boldsymbol{\beta}\|_2^2, \quad (\text{A.19})$$

where we have used that the convex conjugate of $\frac{1}{2} \left(x - \frac{1}{N_0} \right)^2$ is $\frac{1}{2} \left(y + \frac{1}{N_0} \right)^2 - \frac{1}{2N_0^2}$. Solving for α we see that

$$\sum_{W_i=0} \hat{\alpha} + \hat{\boldsymbol{\beta}}' \mathbf{X}_i + 1 = 1$$

Since the lagged outcomes are centered, this implies that

$$\hat{\alpha} = 0$$

Now solving for $\boldsymbol{\beta}$ we see that

$$\mathbf{X}'_0 \left(\mathbf{1} \frac{1}{N_0} + \mathbf{X}_0 \hat{\boldsymbol{\beta}} \right) + \lambda \hat{\boldsymbol{\beta}} = \mathbf{X}_1$$

This implies that

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'_0 \mathbf{X}_0 + \lambda I)^{-1} \mathbf{X}_1$$

Finally, the weights are the ridge weights

$$\hat{\gamma}_i = \frac{1}{N_0} + \mathbf{X}'_1 (\mathbf{X}'_0 \mathbf{X}_0 + \lambda I)^{-1} \mathbf{X}_i = \hat{\gamma}_i^{\text{ridge}}$$

Similarly, the Lagrangian dual to (14) is

$$\min_{\alpha, \boldsymbol{\beta}} \frac{1}{2} \sum_{W_i=0} (\alpha + \boldsymbol{\beta}' \mathbf{X}_i + \hat{\gamma}_i^{\text{scm}})^2 - (\alpha + \boldsymbol{\beta}' \mathbf{X}_1) + \frac{\lambda}{2} \|\boldsymbol{\beta}\|_2^2, \quad (\text{A.20})$$

where we have used that the convex conjugate of $\frac{1}{2} (x - \hat{\gamma}_i^{\text{scm}})^2$ is $\frac{1}{2} (y + \hat{\gamma}_i^{\text{scm}})^2 - \frac{1}{2} \hat{\gamma}_i^{\text{scm}2}$. Solving for α we see that $\hat{\alpha} = 0$. Now solving for $\boldsymbol{\beta}$ we see that

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'_0 \mathbf{X}_0 + \lambda I)^{-1} (\mathbf{X}_1 - \mathbf{X}'_0 \hat{\boldsymbol{\gamma}}^{\text{scm}})$$

Finally, the weights are the ridge ASCM weights

$$\hat{\gamma}_i = \hat{\gamma}_i^{\text{scm}} + (\mathbf{X}_1 - \mathbf{X}'_0 \hat{\boldsymbol{\gamma}}^{\text{scm}})' (\mathbf{X}'_0 \mathbf{X}_0 + \lambda I)^{-1} \mathbf{X}_i = \hat{\gamma}_i^{\text{aug}}$$

□

Proof of Lemma 3. Notice that

$$\begin{aligned} \mathbf{X}_1 - \mathbf{X}'_0 \hat{\boldsymbol{\gamma}}^{\text{aug}} &= (I - \mathbf{X}'_0 \mathbf{X}_0 (\mathbf{X}'_0 \mathbf{X}_0 + N_0 \lambda I)^{-1}) (\mathbf{X}_1 - \mathbf{X}'_0 \hat{\boldsymbol{\gamma}}^{\text{scm}}) \\ &= N_0 \lambda (\mathbf{X}'_0 \mathbf{X}_0 + N_0 \lambda I)^{-1} (\mathbf{X}_1 - \mathbf{X}'_0 \hat{\boldsymbol{\gamma}}^{\text{scm}}) \\ &= \mathbf{V} \text{diag} \left(\frac{\lambda}{d_j^2 + \lambda} \right) \mathbf{V}' (\mathbf{X}_1 - \mathbf{X}'_0 \hat{\boldsymbol{\gamma}}^{\text{scm}}) \end{aligned}$$

So since \mathbf{V} is orthogonal,

$$\|\mathbf{X}_1 - \mathbf{X}'_0 \hat{\boldsymbol{\gamma}}^{\text{aug}}\|_2 = \left\| \text{diag} \left(\frac{\lambda}{d_j^2 + \lambda} \right) (\widetilde{\mathbf{X}}_1 - \widetilde{\mathbf{X}}'_0 \hat{\boldsymbol{\gamma}}^{\text{scm}}) \right\|_2$$

□

Lemma A.2. The ridge augmented SCM weights with hyperparameter λN_0 , $\hat{\boldsymbol{\gamma}}^{\text{aug}}$, satisfy

$$\|\hat{\boldsymbol{\gamma}}^{\text{aug}}\|_2 \leq \|\hat{\boldsymbol{\gamma}}^{\text{scm}}\|_2 + \frac{1}{\sqrt{N_0}} \left\| \text{diag} \left(\frac{d_j}{d_j^2 + \lambda} \right) (\widetilde{\mathbf{X}}_1 - \widetilde{\mathbf{X}}'_0 \hat{\boldsymbol{\gamma}}^{\text{scm}}) \right\|_2, \quad (\text{A.21})$$

with $\widetilde{\mathbf{X}}_i = \mathbf{V}' \mathbf{X}_i$ as defined in Lemma 3.

Proof of Lemma A.2. Notice that using the singular value decomposition and by the triangle in-

equality,

$$\begin{aligned}
\|\hat{\gamma}^{\text{aug}}\|_2 &= \|\hat{\gamma}^{\text{scm}} + \mathbf{X}_0'(\mathbf{X}_0'\mathbf{X}_0 + \lambda I)^{-1}(\mathbf{X}_1 - \mathbf{X}_0'\hat{\gamma}^{\text{scm}})\|_2 \\
&= \left\| \hat{\gamma}^{\text{scm}} + \mathbf{U} \text{diag} \left(\frac{\sqrt{N_0}d_j}{N_0d_j^2 + \lambda N_0} \right) \mathbf{V}'(\mathbf{X}_1 - \mathbf{X}_0'\hat{\gamma}^{\text{scm}}) \right\|_2 \\
&\leq \|\hat{\gamma}^{\text{scm}}\|_2 + \left\| \text{diag} \left(\frac{d_j}{(d_j^2 + \lambda)\sqrt{N_0}} \right) (\tilde{\mathbf{X}}_1 - \tilde{\mathbf{X}}_0'\hat{\gamma}^{\text{scm}}) \right\|_2.
\end{aligned}$$

□

D.2 Proofs for Sections 5, B.1, and B.2

First, consider a model where the post-treatment control potential outcomes at time T are linear in the lagged outcomes and include a unit specific term ξ_i , i.e.

$$Y_{iT}(0) = \boldsymbol{\beta} \cdot \mathbf{X}_i + \xi_i + \varepsilon_{iT}, \quad (\text{A.22})$$

where ε_{it} are independent sub-Gaussian random variables with scale parameter σ , and are ignorable $\mathbb{E}_{\boldsymbol{\varepsilon}_T}[W_i\varepsilon_{iT}] = \mathbb{E}_{\boldsymbol{\varepsilon}_T}[(1 - W_i)\varepsilon_{iT}] = \mathbb{E}_{\boldsymbol{\varepsilon}_T}[\varepsilon_{iT}] = 0$. Below we will put structure on the unit-specific terms ξ_i , first we write a general finite-sample bound.

Proposition A.1. Under model (A.22) with independent sub-Gaussian noise, for weights $\hat{\gamma}$ independent of the post-treatment residuals $(\varepsilon_{1T}, \dots, \varepsilon_{NT})$ and for any $\delta > 0$,

$$Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i Y_{iT} \leq \underbrace{\|\boldsymbol{\beta}\|_2 \left\| \mathbf{X}_1 - \sum_{W_i=0} \hat{\gamma}_i \mathbf{X}_i \right\|_2}_{\text{imbalance in } X} + \underbrace{\left| \xi_1 - \sum_{W_i=0} \hat{\gamma}_i \xi_i \right|}_{\text{approximation error}} + \underbrace{\delta\sigma(1 + \|\hat{\gamma}\|_2)}_{\text{post-treatment noise}}, \quad (\text{A.23})$$

with probability at least $1 - 2e^{-\frac{\delta^2}{2}}$.

Proof. First, note that the estimation error is

$$Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i Y_{iT} = \boldsymbol{\beta} \cdot \left(X_1 - \sum_{W_i=0} \hat{\gamma}_i X_i \right) + \left(\rho_1 - \sum_{W_i=0} \hat{\gamma}_i \xi_i \right) + \left(\varepsilon_{1T} - \sum_{W_i=0} \hat{\gamma}_i \varepsilon_{iT} \right) \quad (\text{A.24})$$

Now since the weights are independent of ε_{iT} , by the ignorability condition (2) and sub-Gaussianity and independence of ε_{iT} , we see that $\varepsilon_{1T} - \sum_{W_i=0} \hat{\gamma}_i \varepsilon_{iT}$ is sub-Gaussian with scale parameter $\sigma\sqrt{1 + \|\hat{\gamma}\|_2^2} \leq \sigma(1 + \|\hat{\gamma}\|_2)$. Therefore we can bound the second term:

$$P \left(\left| \varepsilon_{1T} - \sum_{W_i=0} \hat{\gamma}_i \varepsilon_{iT} \right| \geq \delta\sigma(1 + \|\hat{\gamma}\|_2) \right) \leq 2 \exp \left(-\frac{\delta^2}{2} \right)$$

The result follows from the triangle inequality and the Cauchy-Schwartz inequality. □

Proof of Proposition 1. Note that under the linear model (18), $\xi_i = 0$ for all i . Now from Lemma 3 we have that

$$\|\mathbf{X}_1 - \mathbf{X}'_0 \hat{\gamma}^{\text{aug}}\|_2 = \left\| \text{diag} \left(\frac{\lambda}{d_j^2 + \lambda} \right) (\widetilde{\mathbf{X}}_1 - \widetilde{\mathbf{X}}'_0 \hat{\gamma}^{\text{scm}}) \right\|_2.$$

Plugging this in to Equation (A.23) completes the proof. \square

Proof of Corollary A.1. This is a direct consequence of Proposition A.1 noting that under the linear model (18), $\xi_i = 0$ for all i . \square

Random approximation error We now consider the case where ξ_i are random. We can use Proposition A.1 to further bound the approximation error. In particular, we'll assume that ξ_i are sub-Gaussian random variables with scale parameter ϖ and $\mathbb{E}_\xi[W_i \xi_i] = \mathbb{E}_\xi[(1 - W_i) \xi_i] = \mathbb{E}_\xi[\xi_i] = 0$.

Lemma A.3. If ξ_i are mean-zero sub-Gaussian random variables with scale parameter ϖ , then for weights $\hat{\gamma}$ and any $\delta > 0$ the approximation error satisfies

$$\left| \xi_1 - \sum_{W_i=0} \hat{\gamma}_i \xi_i \right| \leq \delta \varpi + 2 \|\hat{\gamma}\|_1 \varpi \left(\sqrt{\log 2N_0} + \frac{\delta}{2} \right), \quad (\text{A.25})$$

with probability at least $1 - 4e^{-\frac{\delta^2}{2}}$.

Proof of Lemma A.3. From the triangle inequality and Hölder's inequality we see that

$$\left| \xi_1 - \sum_{W_i=0} \hat{\gamma}_i \xi_i \right| \leq |\xi_1| + \|\hat{\gamma}\|_1 \max_{W_i=0} |\xi_i|.$$

Now since the ξ_i are mean-zero sub-Gaussian with scale parameter ϖ , we have that

$$P(|\xi_1| \geq \delta \varpi) \leq 2e^{-\frac{\delta^2}{2}}$$

Next, from the union bound, the maximum of the N_0 sub-Gaussian variables ρ_2, \dots, ρ_N satisfies

$$P\left(\max_{W_i=0} |\xi_i| \geq 2\varpi \sqrt{\log 2N_0} + \delta\right) \leq 2e^{-\frac{\delta^2}{2\varpi^2}}.$$

Setting $\delta = \delta \varpi$ and combining the two probabilities with the union bound gives the result. \square

Lemma A.4. If ξ_i are mean-zero sub-Gaussian random variables with scale parameter ϖ , for the ridge ASCM weights $\hat{\gamma}^{\text{aug}}$ with hyper-parameter $\lambda^{\text{ridge}} = \lambda N_0$ and for any $\delta > 0$ the approximation error satisfies

$$\left| \xi_1 - \sum_{W_i=0} \hat{\gamma}_i \xi_i \right| \leq 2\varpi \left(\sqrt{\log 2N_0} + \frac{\delta}{2} \right) + \underbrace{(1 + \delta)4\varpi \left\| \text{diag} \left(\frac{d_j}{d_j^2 + \lambda} \right) (\widetilde{\mathbf{X}}_1 - \widetilde{\mathbf{X}}'_0 \hat{\gamma}^{\text{scm}}) \right\|_2}_{\text{excess approximation error}}, \quad (\text{A.26})$$

with probability at least $1 - 4e^{-\frac{\delta^2}{2}} - e^{-2(\log 2 + N_0 \log 5)\delta^2}$.

Proof of Lemma A.4. Again from Hölder's inequality we see that

$$\begin{aligned} \left| \xi_1 - \sum_{W_i=0} \hat{\gamma}_i^{\text{aug}} \xi_i \right| &= |\xi_1| + \left| \sum_{W_i=0} (\hat{\gamma}_i^{\text{scm}} + \hat{\gamma}_i^{\text{aug}} - \hat{\gamma}_i^{\text{scm}}) \xi_i \right| \\ &\leq |\xi_1| + \|\hat{\gamma}^{\text{scm}}\|_1 \max_{W_i=0} |\xi_i| + \|\hat{\gamma}^{\text{aug}} - \hat{\gamma}^{\text{scm}}\|_2 \sqrt{\sum_{W_i=0} \xi_i^2}. \end{aligned}$$

We have bounded the first two terms in Lemma A.3, now it suffices to bound the third term. First, from Lemma A.2 we see that

$$\|\hat{\gamma}^{\text{aug}} - \hat{\gamma}^{\text{scm}}\|_2 = \frac{1}{\sqrt{N_0}} \left\| \text{diag} \left(\frac{d_j}{d_j^2 + \lambda} \right) (\tilde{\mathbf{X}}_1 - \tilde{\mathbf{X}}_0' \hat{\gamma}^{\text{scm}}) \right\|_2.$$

Second, via a standard discretization argument (Wainwright, 2018), we can bound the L^2 norm of the vector (ξ_2, \dots, ξ_N) as

$$P \left(\sqrt{\sum_{W_i=0} \xi_i^2} \geq 2\varpi \sqrt{\log 2 + N_0 \log 5} + \delta \right) \leq 2 \exp \left(-\frac{\delta^2}{2\varpi^2} \right).$$

Setting $\delta = 2\delta\varpi \sqrt{\log 2 + N_0 \log 5}$, noting that $\log 2 + N_0 \log 5 < 4N_0$, we have that

$$\|\hat{\gamma}^{\text{aug}} - \hat{\gamma}^{\text{scm}}\|_2 \sqrt{\sum_{W_i=0} \xi_i^2} \leq (1 + \delta)\varpi 4 \left\| \text{diag} \left(\frac{d_j}{d_j^2 + \lambda} \right) (\tilde{\mathbf{X}}_1 - \tilde{\mathbf{X}}_0' \hat{\gamma}^{\text{scm}}) \right\|_2$$

with probability at least $1 - 2e^{-2(\log 2 + N_0 \log 5)\delta^2}$. Since $\|\hat{\gamma}^{\text{scm}}\|_1 = 1$, combining with Lemma A.3 via the union bound gives the result. \square

Theorem A.3. Under model (A.22) with independent sub-Gaussian noise ε_{iT} , and ξ_i mean-zero sub-Gaussian with scale parameter ϖ , for weights $\hat{\gamma}$ independent of the post-treatment outcomes (Y_{1T}, \dots, Y_{NT}) and for any $\delta > 0$,

$$Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i Y_{iT} \leq \underbrace{\|\beta\|_2 \left\| \mathbf{X}_1 - \sum_{W_i=0} \hat{\gamma}_i \mathbf{X}_i \right\|_2}_{\text{imbalance in } X} + \underbrace{\delta\varpi + 2\|\hat{\gamma}\|_1\varpi \left(\sqrt{\log 2N_0} + \frac{\delta}{2} \right)}_{\text{approximation error}} + \underbrace{\delta\sigma(1 + \|\hat{\gamma}\|_2)}_{\text{post-treatment noise}}, \quad (\text{A.27})$$

with probability at least $1 - 6e^{-\frac{\delta^2}{2}}$.

Proof of Theorem A.3. The Theorem directly follows from Proposition A.1 and Lemma A.3, combining the two probabilistic bounds via the union bound. \square

Theorem A.4. Under model (A.22) with independent sub-Gaussian noise ε_{iT} , and ξ_i mean-zero sub-Gaussian with scale parameter ϖ , for any $\delta > 0$, the ridge ASCM weights with hyperparameter $\lambda^{\text{ridge}} = \lambda N_0$ satisfy the bound

$$\begin{aligned}
Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i Y_{iT} \leq & \underbrace{\|\beta\|_2 \left\| \text{diag} \left(\frac{\lambda}{d_j^2 + \lambda} \right) \left(\widetilde{\mathbf{X}}_1 - \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} \widetilde{\mathbf{X}}_i \right) \right\|_2}_{\text{imbalance in } X} + \underbrace{2\varpi \left(\sqrt{\log 2N_0} + \frac{\delta}{2} \right)}_{\text{approximation error}} \\
& + \underbrace{(1 + \delta)4\varpi \left\| \text{diag} \left(\frac{d_j}{d_j^2 + \lambda} \right) \left(\widetilde{\mathbf{X}}_1 - \widetilde{\mathbf{X}}_0' \hat{\gamma}^{\text{scm}} \right) \right\|_2}_{\text{excess approximation error}} + \underbrace{\delta\sigma (1 + \|\hat{\gamma}\|_2)}_{\text{post-treatment noise}},
\end{aligned} \tag{A.28}$$

with probability at least $1 - 6e^{-\frac{\delta^2}{2}} - e^{-2(\log 2 + N_0 \log 5)\delta^2}$.

Proof of Theorem A.4. First note that from Lemma 3 we have that

$$\|\mathbf{X}_1 - \mathbf{X}_0' \hat{\gamma}^{\text{aug}}\|_2 = \left\| \text{diag} \left(\frac{\lambda}{d_j^2 + \lambda} \right) (\widetilde{\mathbf{X}}_1 - \widetilde{\mathbf{X}}_0' \hat{\gamma}^{\text{scm}}) \right\|_2.$$

The Theorem directly follows from Proposition A.1 and Lemma A.4, combining the two probabilistic bounds via the union bound. \square

Theorems A.3 and A.4 have several implications when the outcomes follow a linear factor model (21). To see this, we need one more lemma:

Lemma A.5. The linear factor model is a special case of model (A.22) with $\beta = \frac{1}{T_0} \boldsymbol{\mu} \boldsymbol{\mu}'$ and $\xi_i = \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu} \varepsilon_{i(1:T_0)}$. $\|\beta\|_2 \leq \frac{MJ^2}{\sqrt{T_0}}$, and if $\varepsilon_{i(1:T_0)}$ are independent sub-Gaussian vectors with scale parameter σ_{T_0} , then $\frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' \varepsilon_{i(1:T_0)}$ is sub-Gaussian with scale parameter $\frac{JM^2 \sigma_{T_0}}{\sqrt{T_0}}$.

Proof of Lemma A.5. Notice that under the linear factor model, the pre-treatment covariates for unit i satisfy:

$$\mathbf{X}_i = \boldsymbol{\mu} \phi_i + \varepsilon_{i(1:T_0)}.$$

Multiplying both sides by $(\boldsymbol{\mu}' \boldsymbol{\mu})^{-1} \boldsymbol{\mu}' = \frac{1}{T_0} \boldsymbol{\mu}'$ and rearranging gives

$$\frac{1}{T_0} \boldsymbol{\mu}' \mathbf{X}_i - \frac{1}{T_0} \boldsymbol{\mu}' \varepsilon_{i(1:T_0)} = \phi_i.$$

Then we see that the post treatment outcomes are

$$Y_{iT}(0) = \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' \mathbf{X}_i + \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' \varepsilon_{i(1:T_0)}.$$

Since $\varepsilon_{i(1:T_0)}$ is a sub-Gaussian vector $v' \varepsilon_{i(1:T_0)}$ is sub-Gaussian with scale parameter σ_{T_0} for any $v \in \mathbb{R}^{T_0}$ such that $\|v\|_2 = 1$. Now notice that $\|\boldsymbol{\mu}'_T \boldsymbol{\mu}'\|_2 \leq \|\boldsymbol{\mu}_T\|_2 \|\boldsymbol{\mu}\|_2 \leq MJ^2 \sqrt{T_0}$. This completes the proof. \square

Proof of Corollary A.2. From Lemma A.5 we can apply Theorem A.3 with $\beta = \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}'$ and $\xi_i = \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' \boldsymbol{\varepsilon}_{i(1:T_0)}$. Since ε_{it} are independent sub-Gaussian random variables, $\boldsymbol{\varepsilon}_{i(1:T_0)}$ is a sub-Gaussian vector with scale parameter $\sigma_{T_0} = \sigma$. Noting that $\|\hat{\boldsymbol{\gamma}}\|_1 = \sum_{W_i=0} |\hat{\gamma}_i| = 1$ and applying Lemma A.5 gives the result. \square

Proof of Theorem 1. Again from Lemma A.5 we can apply Theorem A.4 with $\beta = \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}'$ and $\xi_i = \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' \boldsymbol{\varepsilon}_{i(1:T_0)}$, so $\varpi = \frac{JM^2}{\sqrt{T_0}}$. Plugging these values into Theorem A.3 gives the result. \square

Corollary A.5 (Approximation error for ridge ASCM with dependent errors). Under the linear factor model (21) with time-dependent errors satisfying $\boldsymbol{\varepsilon}_{i(1:T_0)} \stackrel{iid}{\sim} N(0, \sigma^2 \Omega)$ the approximation error satisfies

$$\begin{aligned} \left| \xi_1 - \sum_{W_i=0} \hat{\gamma}_i \xi_i \right| &= \left| \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' \left(\boldsymbol{\varepsilon}_{1(1:T_0)} - \sum_{W_i=0} \hat{\gamma}_i \boldsymbol{\varepsilon}_{i(1:T_0)} \right) \right| \\ &\leq 2 \sqrt{\frac{\|\Omega\|_2}{T_0}} JM^2 \sigma \left(\sqrt{\log 2N_0} + \delta + (1 + \delta) 2 \left\| \text{diag} \left(\frac{d_j}{d_j^2 + \lambda} \right) \left(\widetilde{\mathbf{X}}_1 - \widetilde{\mathbf{X}}_0 \cdot \hat{\boldsymbol{\gamma}}^{\text{scm}} \right) \right\|_2 \right), \end{aligned} \quad (\text{A.29})$$

Proof of Corollary A.5. From Lemma A.5, we see that $\xi_i = \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' \boldsymbol{\varepsilon}_{i(1:T_0)}$ is sub-Gaussian with scale parameter $JM^2 \sqrt{\frac{\|\Omega\|_2}{T_0}}$. Plugging in to Lemma A.4 gives the result. \square

Lipshitz approximation error If ξ_i are Lipshitz functions, we can also bound the approximation error. Specifically, we assume that $\xi_i = f(\mathbf{X}_i)$ where $f : \mathbb{R}^{T_0} \rightarrow \mathbb{R}$ is an L -Lipshitz function.

Lemma A.6. If $\xi_i = f(\mathbf{X}_i)$ where $f : \mathbb{R}^{T_0} \rightarrow \mathbb{R}$ is an L -Lipshitz function, then for weights on the simplex $\hat{\boldsymbol{\gamma}} \in \Delta^{N_0}$, the approximation error satisfies

$$\left| \xi_1 - \sum_{W_i=0} \hat{\gamma}_i \xi_i \right| \leq L \sum_{W_i=0} \hat{\gamma}_i \|\mathbf{X}_1 - \mathbf{X}_i\|_2 \quad (\text{A.30})$$

Proof of Lemma A.6. Since the weights sum to one, we have that

$$\left| \xi_1 - \sum_{W_i=0} \hat{\gamma}_i \xi_i \right| \leq \left| \sum_{W_i=0} \hat{\gamma}_i (f(\mathbf{X}_1) - f(\mathbf{X}_i)) \right|.$$

Now from the Lipshitz property, $|f(\mathbf{X}_1) - f(\mathbf{X}_i)| \leq L \|\mathbf{X}_1 - \mathbf{X}_i\|_2$, and so by Jensen's inequality:

$$\left| \sum_{W_i=0} \hat{\gamma}_i (f(\mathbf{X}_1) - f(\mathbf{X}_i)) \right| \leq L \sum_{W_i=0} \hat{\gamma}_i \|\mathbf{X}_1 - \mathbf{X}_i\|_2$$

\square

Proof of Theorem A.3. The proof follows directly from Proposition A.1 and Lemma A.6. \square

Lemma A.7. Let $C = \max_{W_i=0} \|\mathbf{X}_i\|_2$. If $\xi_i = f(\mathbf{X}_i)$ where $f : \mathbb{R}^{T_0} \rightarrow \mathbb{R}$ is an L -Lipshitz function, then the ridge ASCM weights $\hat{\gamma}^{\text{aug}}$ (13) with hyperparameter $\lambda^{\text{ridge}} = N_0\lambda$ satisfy

$$\left| \xi_1 - \sum_{W_i=0} \hat{\gamma}_i^{\text{aug}} \xi_i \right| \leq L \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} \|\mathbf{X}_1 - \mathbf{X}_i\|_2 + CL \left\| \text{diag} \left(\frac{d_j}{d_j^2 + \lambda} \right) (\tilde{\mathbf{X}}_1 - \tilde{\mathbf{X}}_0' \hat{\gamma}^{\text{scm}}) \right\|_2 \quad (\text{A.31})$$

Proof of Lemma A.7. From the triangle inequality we have that

$$\left| \xi_1 - \sum_{W_i=0} \hat{\gamma}_i^{\text{aug}} \xi_i \right| \leq \left| \sum_{W_i=0} \hat{\gamma}_i^{\text{scm}} (f(\mathbf{X}_1) - f(\mathbf{X}_i)) \right| + \left| \sum_{W_i=0} \mathbf{X}_i (\mathbf{X}_0' \mathbf{X}_0 + \lambda I)^{-1} (\mathbf{X}_1 - \mathbf{X}_0' \hat{\gamma}^{\text{scm}}) f(\mathbf{X}_i) \right|.$$

We have already bounded the first term in Lemma A.6, now we bound the second term. From the Cauchy-Schwartz inequality and since $\|x\|_2 \leq \sqrt{N_0} \|x\|_\infty$ for all $x \in \mathbb{R}^{N_0}$ we see that

$$\begin{aligned} \left| \sum_{W_i=0} \mathbf{X}_i (\mathbf{X}_0' \mathbf{X}_0 + \lambda I)^{-1} (\mathbf{X}_1 - \mathbf{X}_0' \hat{\gamma}^{\text{scm}}) f(\mathbf{X}_i) \right| &\leq \sqrt{N_0} \left\| \mathbf{X}_0 (\mathbf{X}_0' \mathbf{X}_0 + \lambda I)^{-1} (\mathbf{X}_1 - \mathbf{X}_0' \hat{\gamma}^{\text{scm}}) \right\|_2 |f(\mathbf{X}_i)| \\ &= \left\| \text{diag} \left(\frac{d_j}{d_j^2 + \lambda} \right) (\tilde{\mathbf{X}}_1 - \tilde{\mathbf{X}}_0' \hat{\gamma}^{\text{scm}}) \right\|_2 |f(\mathbf{X}_i)| \\ &\leq CL \left\| \text{diag} \left(\frac{d_j}{d_j^2 + \lambda} \right) (\tilde{\mathbf{X}}_1 - \tilde{\mathbf{X}}_0' \hat{\gamma}^{\text{scm}}) \right\|_2, \end{aligned}$$

where the second line comes from Lemma A.2 and the third line from the Lipshitz property. \square

Proof of Theorem A.1. The proof follows directly from Proposition A.1 and Lemma A.7. \square

D.3 Proofs for Sections 6 and B.3

Proof of Lemma 4. The regression parameters $\hat{\eta}_x$ and $\hat{\eta}_z$ in Equation (27) are:

$$\hat{\eta}_x = (\check{\mathbf{X}}_0' \check{\mathbf{X}}_0 + \lambda^{\text{ridge}} I)^{-1} \check{\mathbf{X}}_0' Y_{0T} \quad \text{and} \quad \hat{\eta}_z = (\mathbf{Z}_0' \mathbf{Z}_0)^{-1} \mathbf{Z}_0' Y_{0T} \quad (\text{A.32})$$

Now notice that

$$\begin{aligned} \hat{Y}_{0T}^{\text{cov}} &= \hat{\eta}_x' \mathbf{X}_1 + \hat{\eta}_z' \mathbf{Z}_1 + \sum_{W_i=0} (Y_{iT} - \hat{\eta}_x' \mathbf{X}_i - \hat{\eta}_z' \mathbf{Z}_i) \hat{\gamma}_i \\ &= \hat{\eta}_x' (\mathbf{X}_1 - \mathbf{X}_0' \hat{\gamma}) + \hat{\eta}_z' (\mathbf{Z}_1 - \mathbf{Z}_0' \hat{\gamma}) + \mathbf{Y}_{0T}' \hat{\gamma} \\ &= \hat{\eta}_x' (\mathbf{X}_1 - \mathbf{X}_0' \hat{\gamma}) - \hat{\eta}_x' \mathbf{X}_0 (\mathbf{Z}_0' \mathbf{Z}_0)^{-1} (\mathbf{Z}_1 - \mathbf{Z}_0' \hat{\gamma}) + Y_{0T}' \mathbf{Z}_0 (\mathbf{Z}_0' \mathbf{Z}_0)^{-1} (\mathbf{Z}_1 - \mathbf{Z}_0' \hat{\gamma}) + Y_{0T}' \hat{\gamma} \\ &= \hat{\eta}_x' (\check{\mathbf{X}}_1 - \check{\mathbf{X}}_0' \hat{\gamma}) + Y_{0T}' \mathbf{Z}_0 (\mathbf{Z}_0' \mathbf{Z}_0)^{-1} (\mathbf{Z}_1 - \mathbf{Z}_0' \hat{\gamma}) + Y_{0T}' \hat{\gamma} \\ &= Y_{0T}' \left(\hat{\gamma} + \check{\mathbf{X}}_0 (\check{\mathbf{X}}_0' \check{\mathbf{X}}_0 + \lambda^{\text{ridge}} I_{T_0})^{-1} (\check{\mathbf{X}}_1 - \check{\mathbf{X}}_0' \hat{\gamma}) + \mathbf{Z}_0 (\mathbf{Z}_0' \mathbf{Z}_0)^{-1} (\mathbf{Z}_1 - \mathbf{Z}_0' \hat{\gamma}) \right). \end{aligned} \quad (\text{A.33})$$

This gives the form of $\hat{\gamma}^{\text{cov}}$. The imbalance in Z is

$$\begin{aligned} \mathbf{Z}_1 - \mathbf{Z}'_0 \hat{\gamma}^{\text{cov}} &= (\mathbf{Z}_1 - \mathbf{Z}'_0 \mathbf{Z}_0 (\mathbf{Z}'_0 \mathbf{Z}_0)^{-1} \mathbf{Z}_1) + (\mathbf{Z}_0 - \mathbf{Z}'_0 \mathbf{Z}_0 (\mathbf{Z}'_0 \mathbf{Z}_0)^{-1} \mathbf{Z}_0)' \hat{\gamma} \\ &\quad - \mathbf{Z}'_0 \check{\mathbf{X}}_0 (\check{\mathbf{X}}'_0 \check{\mathbf{X}}_0 + \lambda^{\text{ridge}} \mathbf{I})^{-1} (\check{\mathbf{X}}_1 - \check{\mathbf{X}}'_0 \hat{\gamma}) \\ &= 0. \end{aligned} \tag{A.34}$$

The pre-treatment fit is

$$\begin{aligned} \mathbf{X}_1 - \mathbf{X}'_0 \hat{\gamma}^{\text{cov}} &= (\mathbf{X}_1 - \mathbf{X}'_0 \mathbf{Z}_0 (\mathbf{Z}'_0 \mathbf{Z}_0)^{-1} \mathbf{Z}_1) + (\mathbf{X}_0 - \mathbf{X}'_0 \mathbf{Z}_0 (\mathbf{Z}'_0 \mathbf{Z}_0)^{-1} \mathbf{Z}_0)' \hat{\gamma} \\ &\quad - \mathbf{X}'_0 \check{\mathbf{X}}_0 (\check{\mathbf{X}}'_0 \check{\mathbf{X}}_0 + \lambda^{\text{ridge}} \mathbf{I}_{T_0})^{-1} (\check{\mathbf{X}}_1 - \check{\mathbf{X}}'_0 \hat{\gamma}) \\ &= \left(\mathbf{I}_{T_0} - \mathbf{X}'_0 \check{\mathbf{X}}_0 (\check{\mathbf{X}}'_0 \check{\mathbf{X}}_0 + \lambda^{\text{ridge}} \mathbf{I}_{T_0})^{-1} \right) (\check{\mathbf{X}}_1 - \check{\mathbf{X}}'_0 \hat{\gamma}) \\ &= \left(\mathbf{I}_{T_0} - \check{\mathbf{X}}'_0 \check{\mathbf{X}}_0 (\check{\mathbf{X}}'_0 \check{\mathbf{X}}_0 + \lambda^{\text{ridge}} \mathbf{I}_{T_0})^{-1} \right) (\check{\mathbf{X}}_1 - \check{\mathbf{X}}'_0 \hat{\gamma}). \end{aligned} \tag{A.35}$$

This gives the bound on the pre-treatment fit. □

Proof of Theorem A.2. First, we will separate $f(\mathbf{Z})$ into the projection onto \mathbf{Z} and a residual. Defining $\mathbf{B}_t = (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}' f_t(\mathbf{Z}) \in \mathbb{R}^K$ as the regression coefficient, the projection of $f_t(\mathbf{Z}_i)$ is $\mathbf{Z}'_i \mathbf{B}_t$ and the residual is $e_{it} = f_t(\mathbf{Z}_i) - \mathbf{Z}'_i \mathbf{B}_t$. We will denote the matrix of regression coefficients over $t = 1, \dots, T_0$ as $\mathbf{B} = [\mathbf{B}_1, \dots, \mathbf{B}_{T_0}] \in \mathbb{R}^{K \times T_0}$ and denote the matrix of residuals as $\mathbf{E} \in \mathbb{R}^{n \times T_0}$, with $\mathbf{E}_1 = (e_{11}, \dots, e_{1T_0})$ as the vector of residuals for the treated unit and \mathbf{E}_0 as the matrix of residuals for the control units.

Then the error is

$$\begin{aligned} \left| Y_{1T}(0) - \sum_{W_i=0} \hat{\gamma}_i^{\text{cov}} Y_{iT} \right| &\leq \left| \boldsymbol{\mu}_T \cdot \left(\phi_1 - \sum_{W_i=0} \hat{\gamma}_i^{\text{cov}} \phi_i \right) \right| + \left| \mathbf{B}_t \cdot \left(\mathbf{Z}_1 - \sum_{W_i=0} \hat{\gamma}_i^{\text{cov}} \mathbf{Z}_i \right) \right| \\ &\quad + \left| e_{1T} - \sum_{W_i=0} \hat{\gamma}_i^{\text{cov}} e_{iT} \right| + \left| \varepsilon_{1T} - \sum_{W_i=0} \hat{\gamma}_i^{\text{cov}} \varepsilon_{iT} \right| \end{aligned}$$

Since $\hat{\gamma}_i^{\text{cov}}$ exactly balances the covariates, the second term is equal to zero. We can bound the third term with Hölder's inequality:

$$\left| e_{1T} - \sum_{W_i=0} \hat{\gamma}_i^{\text{cov}} e_{iT} \right| \leq |e_{1T}| + \sqrt{RSS_T} \|\hat{\gamma}^{\text{cov}}\|_2$$

In previous theorems we have bounded the last term with high probability. Only the error due to imbalance remains.

Denote $\boldsymbol{\varepsilon}_{0(1:T_0)}$ as the matrix of pre-treatment noise for the control units, where the rows correspond to $\boldsymbol{\varepsilon}_{2(1:T_0)}, \dots, \boldsymbol{\varepsilon}_{N_0(1:T_0)}$. Building on Lemma A.5, we can see that the error due to

imbalance in ϕ is equal to

$$\begin{aligned} \boldsymbol{\mu}_T \cdot \left(\phi_1 - \sum_{W_i=0} \hat{\gamma}_i^{\text{cov}} \phi_i \right) &= \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' (\mathbf{X}_1 - \mathbf{X}'_0 \hat{\gamma}^{\text{cov}}) - \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' (\boldsymbol{\varepsilon}_{1(1:T_0)} - \boldsymbol{\varepsilon}'_{0(1:T_0)} \hat{\gamma}^{\text{cov}}) \\ &\quad - \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' B' (\mathbf{Z}_1 - \mathbf{Z}'_0 \hat{\gamma}^{\text{cov}}) - \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' (\mathbf{E}_1 - \mathbf{E}'_0 \hat{\gamma}^{\text{cov}}). \end{aligned} \quad (\text{A.36})$$

By construction, $\hat{\gamma}^{\text{cov}}$ perfectly balances the covariates, and combined with Lemma 4, the error due to imbalance in ϕ simplifies to

$$\boldsymbol{\mu}_T \cdot \left(\phi_1 - \sum_{W_i=0} \gamma_i \phi_i \right) = \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' (\tilde{\mathbf{X}}_1 - \tilde{\mathbf{X}}'_0 \hat{\gamma}) - \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' (\boldsymbol{\varepsilon}_{1(1:T_0)} - \boldsymbol{\varepsilon}'_{0(1:T_0)} \hat{\gamma}^{\text{cov}}) - \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' (\mathbf{E}_1 - \mathbf{E}'_0 \hat{\gamma}^{\text{cov}}).$$

We now turn to bounding the noise term and the error due to the projection of $f(Z)$ on to Z . First, notice that

$$\frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' \boldsymbol{\varepsilon}'_{0(1:T_0)} \hat{\gamma}^{\text{cov}} = \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' \boldsymbol{\varepsilon}'_{0(1:T_0)} \hat{\gamma}^{\text{scm}} + \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' \boldsymbol{\varepsilon}'_{0(1:T_0)} \mathbf{Z}_0 (\mathbf{Z}'_0 \mathbf{Z}_0)^{-1} (\mathbf{Z}_1 - \mathbf{Z}'_0 \hat{\gamma}^{\text{scm}}).$$

We have bounded the first term on the right hand side in Lemma A.3. To bound the second term, notice that $\sum_{W_i=0} \sum_{t=1}^{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}_t Z_{ik} \varepsilon_{it}$ is sub-Gaussian with scale parameter $\sigma M J^2 \sqrt{T_0} \|Z_{\cdot k}\|_2 = M J^2 \sigma \sqrt{T_0 N_0}$. We can now bound the L^2 norm of $\frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' \boldsymbol{\varepsilon}'_{0(1:T_0)} \mathbf{Z}_0 \in \mathbb{R}^K$:

$$P \left(\frac{1}{T_0} \|\boldsymbol{\mu}'_T \boldsymbol{\mu}' \boldsymbol{\varepsilon}'_{0(1:T_0)} \mathbf{Z}_0\|_2 \geq 2JM^2\sigma \left(\sqrt{\frac{N_0 K \log 5}{T_0}} + \delta \right) \right) \leq 2 \exp \left(-\frac{T_0 \delta^2}{2} \right)$$

Replacing δ with $\sqrt{\frac{KN_0}{T_0}}(2 - \sqrt{\log 5})$ and with the Cauchy-Schwarz inequality we see that

$$\frac{1}{T_0} \left| \boldsymbol{\mu}'_T \boldsymbol{\mu}' \boldsymbol{\varepsilon}'_{0(1:T_0)} \mathbf{Z}_0 (\mathbf{Z}'_0 \mathbf{Z}_0)^{-1} (\mathbf{Z}_1 - \mathbf{Z}'_0 \hat{\gamma}) \right| \leq 4JM^2\sigma \sqrt{\frac{K}{T_0 N_0}} \|\mathbf{Z}_1 - \mathbf{Z}'_0 \hat{\gamma}^{\text{scm}}\|_2$$

with probability at least $1 - 2 \exp \left(-\frac{KN_0(2-\sqrt{\log 5})^2}{2} \right)$.

Next we turn to the residual term. By Hölder's inequality and using that for a matrix \mathbf{A} , the operator norm is bounded by $\|\mathbf{A}\|_2 \leq \sqrt{\text{trace}(\mathbf{A}'\mathbf{A})}$ we see that

$$\begin{aligned} \left| \frac{1}{T_0} \boldsymbol{\mu}'_T \boldsymbol{\mu}' (\mathbf{E}_1 - \mathbf{E}'_0 \hat{\gamma}^{\text{cov}}) \right| &\leq \frac{JM^2}{\sqrt{T_0}} (\|\mathbf{E}_1\|_2 + \|\hat{\gamma}^{\text{cov}}\|_2 \|\mathbf{E}_0\|_2) \\ &\leq JM^2 \left(\max_{t=1, \dots, T_0} |e_{1t}| + \|\hat{\gamma}^{\text{cov}}\|_2 \sqrt{\frac{1}{T_0} \sum_{t=1}^{T_0} RSS_t} \right) \\ &\leq JM^2 \left(\max_{t=1, \dots, T_0} |e_{1t}| + \|\hat{\gamma}^{\text{cov}}\|_2 \sqrt{\max_t RSS_t} \right), \end{aligned}$$

where we have used that $\frac{1}{\sqrt{T_0}} \|\mathbf{E}_1\|_2 \leq \max_{t=1, \dots, T_0} |e_{1t}|$ and $\text{trace}(\mathbf{E}'_0 \mathbf{E}_0) = \sum_{t=1}^{T_0} RSS_t$.

Combining with Lemma 4 and putting together the pieces with the union bound gives the result. \square

E Connection to balancing weights and IPW

We have motivated Augmented SCM via bias correction. An alternative motivation comes from the connection between SCM and inverse propensity score weighting (IPW).

First, notice that the SCM weights from the constrained optimization problem in Equation (4) are a form of *approximate balancing weights* (see, for example Zubizarreta, 2015; Athey et al., 2018; Tan, 2017; Wang and Zubizarreta, 2018; Zhao, 2018). Unlike traditional inverse propensity score weights, which indirectly minimize covariate imbalance by estimating a propensity score model, balancing weights seek to *directly* minimize covariate imbalance, in this case L^2 imbalance. Balancing weights have a Lagrangian dual formulation as inverse propensity score weights (see, for example Zhao and Percival, 2017; Zhao, 2018; Chattopadhyay et al., 2020). Extending these results to the SCM setting, the Lagrangian dual of the SCM optimization problem in Equation (4) has the form of a propensity score model. Importantly, as we discuss below, it is not always appropriate to interpret this model as a propensity score.

We first derive the Lagrangian dual for a general class of balancing weights problems, then specialize to the penalized SCM estimator (4).

$$\begin{aligned} \min_{\boldsymbol{\gamma}} \quad & \underbrace{h_{\zeta}(\mathbf{X}_1 - \mathbf{X}'_0 \boldsymbol{\gamma})}_{\text{balance criterion}} + \sum_{W_i=0} \underbrace{f(\gamma_i)}_{\text{dispersion}} \\ \text{subject to} \quad & \sum_{W_i=0} \gamma_i = 1. \end{aligned} \tag{A.37}$$

This formulation generalizes Equation (4) in two ways: first, we remove the non-negativity constraint and note that this can be included by restricting the domain of the strongly convex dispersion penalty f . Examples include the re-centered L^2 dispersion penalties for ridge regression and ridge ASCM, an entropy penalty (Robbins et al., 2017), and an elastic net penalty (Doudchenko and Imbens, 2017). Second, we generalize from the squared L^2 norm to a general balance criterion h_{ζ} ; another prominent example is an L^{∞} constraint (see e.g. Zubizarreta, 2015; Athey et al., 2018).

Proposition A.2. The Lagrangian dual to Equation (A.37) is

$$\min_{\alpha, \boldsymbol{\beta}} \quad \underbrace{\sum_{W_i=0} f^*(\alpha + \boldsymbol{\beta}' X_i) - (\alpha + \boldsymbol{\beta}' \mathbf{X}_1)}_{\text{loss function}} + \underbrace{h_{\zeta}^*(\boldsymbol{\beta})}_{\text{regularization}}, \tag{A.38}$$

where a convex, differentiable function g has convex conjugate $g^*(\mathbf{y}) \equiv \sup_{\mathbf{x} \in \text{dom}(g)} \{\mathbf{y}' \mathbf{x} - g(\mathbf{x})\}$. The solutions to the primal problem (A.37) are $\hat{\gamma}_i = f^{*'}(\hat{\alpha} + \hat{\boldsymbol{\beta}}' X_i)$, where $f^{*'}(\cdot)$ is the first derivative of the convex conjugate, $f^*(\cdot)$.

There is a large literature relating balancing weights to propensity score weights. This literature shows that the loss function in Equation (A.38) is an M-estimator for the propensity score and thus will be consistent for the propensity score parameters under large N asymptotics. The dispersion measure $f(\cdot)$ determines the link function of the propensity score model, where the odds of treatment are $\frac{\pi(x)}{1-\pi(x)} = f^{*'}(\alpha + \boldsymbol{\beta}' x)$. Note that un-penalized SCM, which can yield multiple solutions, does not have a well-defined link function. We extend the duality to a general set of balance criteria so that Equation (A.38) is a regularized M-estimator of the propensity score parameters

where the balance criterion $h_\zeta(\cdot)$ determines the type of regularization through its conjugate $h_\zeta^*(\cdot)$. This formulation recovers the duality between entropy balancing and a logistic link (Zhao and Percival, 2017), Oaxaca-Blinder weights and a log-logistic link (Kline, 2011), and L^∞ balance and L^1 regularization (Wang and Zubizarreta, 2018). This more general formulation also suggests natural extensions of both SCM and ASCM beyond the L^2 setting to other forms, especially L^1 regularization.

Specializing proposition A.2 to a squared L^2 balance criterion $h_\zeta(x) = \frac{1}{2\zeta} \|x\|_2^2$ as in the penalized SCM problems yields that the dual propensity score coefficients β are regularized by a ridge penalty. In the case of an entropy dispersion penalty as Robbins et al. (2017) consider, the donor weights $\hat{\gamma}$ have the form of IPW weights with a logistic link function, where the propensity score is $\pi(\mathbf{X}_i) = \text{logit}^{-1}(\alpha + \beta' \mathbf{X}_i)$, the odds of treatment are $\frac{\pi(\mathbf{X}_i)}{1-\pi(\mathbf{X}_i)} = \exp(\alpha + \beta' \mathbf{X}_i) = \gamma_i$.

We emphasize that while Proposition A.2 shows that the the estimated weights have the IPW form, in SCM settings it may not always be appropriate to interpret the dual problem as a propensity score reflecting stochastic selection into treatment. For example, this interpretation would not be appropriate in some canonical SCM examples, such as the analysis of German reunification in Abadie et al. (2015).

Proof of Proposition A.2. We can augment the optimization problem (A.37) with auxiliary variables ϵ , yielding:

$$\begin{aligned} \min_{\gamma, \epsilon} \quad & h_\zeta(\epsilon) + \sum_{W_i=0} f(\gamma_i). \\ \text{subject to } \quad & \epsilon = \mathbf{X}_1 - \mathbf{X}'_0 \gamma \\ & \sum_{W_i=0} \gamma_i = 1 \end{aligned} \tag{A.39}$$

The Lagrangian is

$$\mathcal{L}(\gamma, \epsilon, \alpha, \beta) = \sum_{i|W_i=0} f(\gamma_i) + \alpha(1 - \gamma_i) + h_\zeta(\epsilon) + \beta'(\mathbf{X}_1 - \mathbf{X}'_0 \gamma - \epsilon). \tag{A.40}$$

The dual maximizes the objective

$$\begin{aligned} q(\alpha, \beta) &= \min_{\gamma, \epsilon} \mathcal{L}(\gamma, \epsilon, \alpha, \beta) \\ &= \sum_{W_i=0} \min_{\gamma_i} \{f(\gamma_i) - (\alpha + \beta' \mathbf{X}_i) \gamma_i\} + \min_{\epsilon} \{h_\zeta(\epsilon) - \beta' \epsilon\} + \alpha + \beta' \mathbf{X}_1 \\ &= - \sum_{W_i=0} f^*(\alpha + \beta' \mathbf{X}_i) + \alpha + \beta' \mathbf{X}'_1 - h_\zeta^*(\beta), \end{aligned} \tag{A.41}$$

By strong duality the general dual problem (A.38), which minimizes $-q(\alpha, \beta)$, is equivalent to the primal balancing weights problem. Given the $\hat{\alpha}$ and $\hat{\beta}$ that minimize the Lagrangian dual objective, $-q(\alpha, \beta)$, we recover the donor weights solution to (A.37) as

$$\hat{\gamma}_i = f^{*'}(\hat{\alpha} + \hat{\beta}' \mathbf{X}_i). \tag{A.42}$$

□

F Additional figures

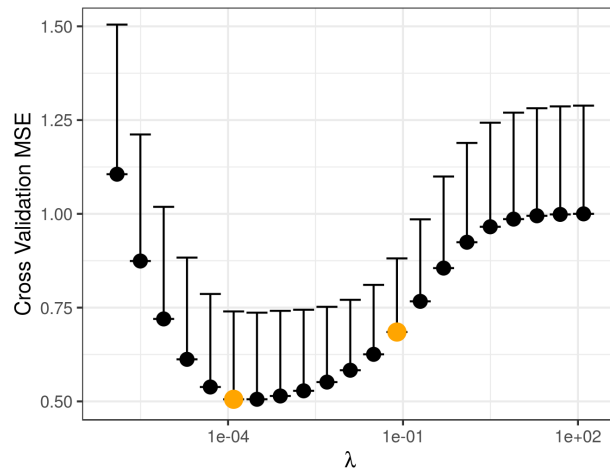


Figure F.1: Cross validation MSE and one standard error computed according to Equation (25). The minimal point, and the maximum λ within one standard error of the minimum are highlighted.

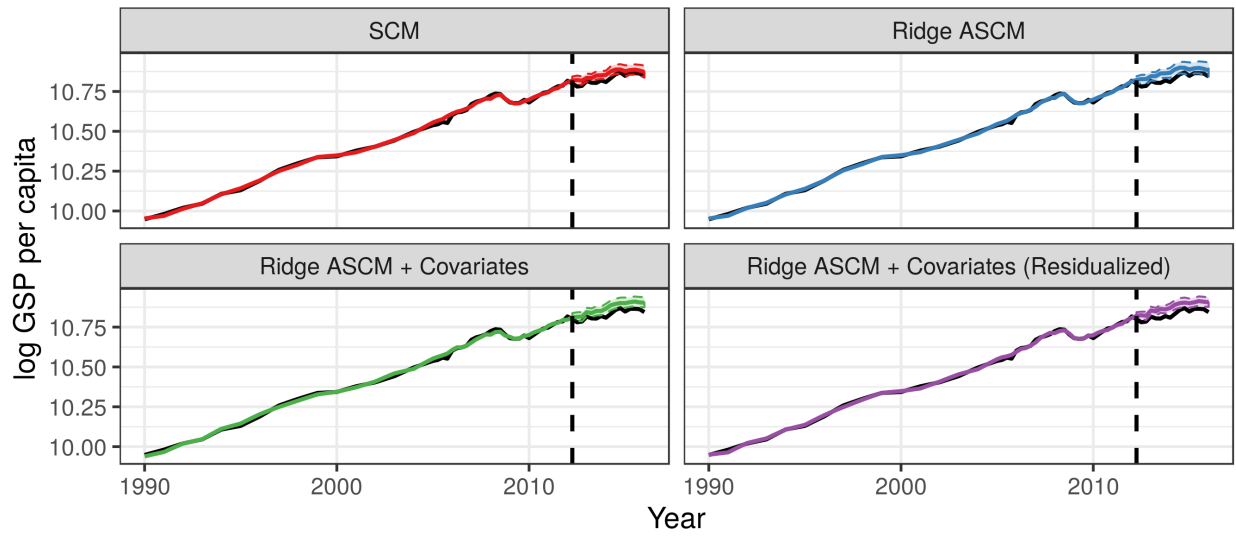


Figure F.2: Point estimates along with point-wise 95% conformal prediction intervals for counterfactual log GSP per capita without the tax cuts using SCM, ridge ASCM, and ridge ASCM with covariates, plotting with the observed log GSP per capita in black.

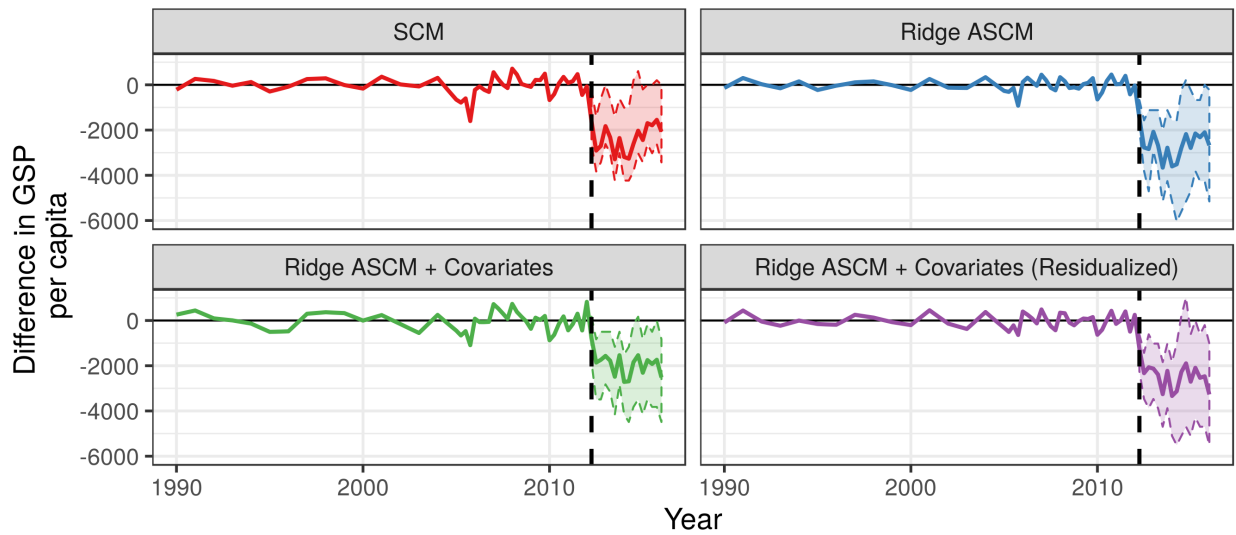


Figure F.3: Point estimates along with point-wise 95% conformal confidence intervals for the effect of the tax cuts on GSP per capita using SCM, ridge ASCM, and ridge ASCM with covariates.

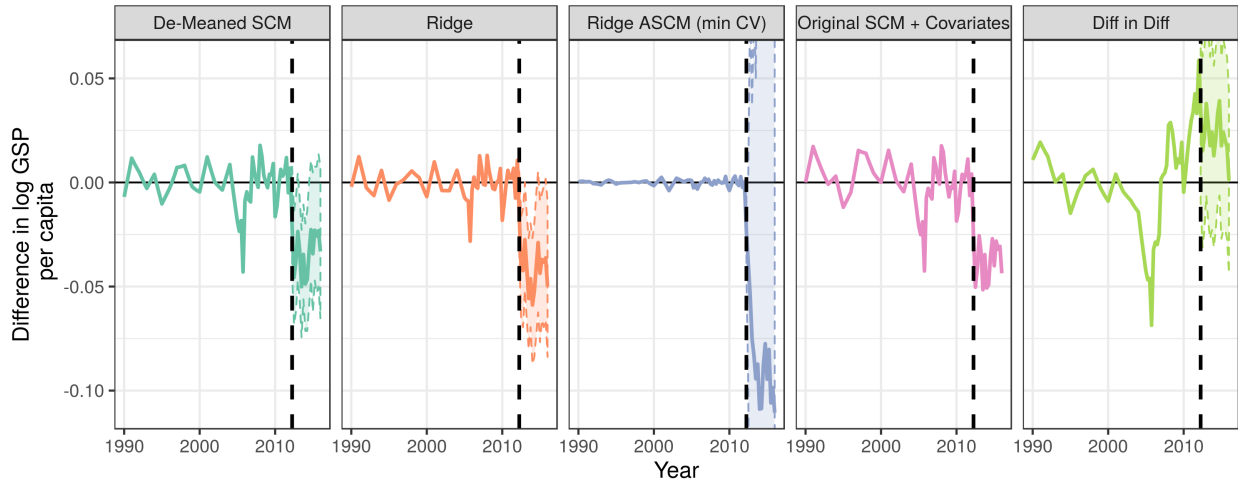


Figure F.4: Point estimates \pm two standard errors of the ATT for the effect of the tax cuts on log GSP per capita using de-meaned SCM, ridge regression alone, ridge ASCM with λ chosen to minimize the cross validated MSE, the original SCM proposal with covariates (Abadie et al., 2010), and a two-way fixed effects differences in differences estimate.

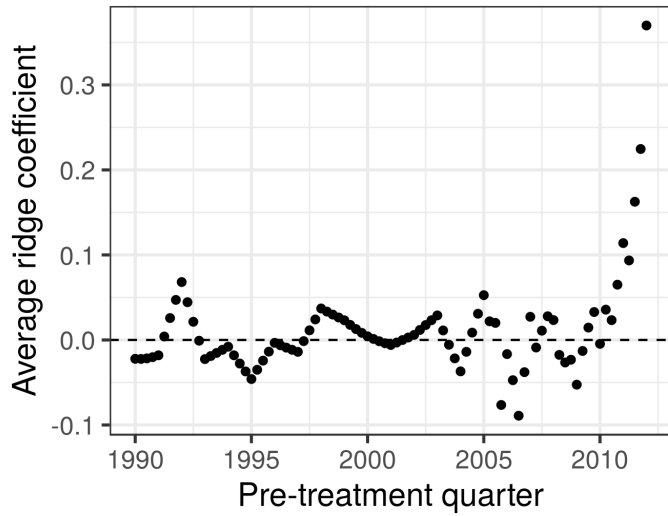


Figure F.5: Ridge regression coefficients for each pre-treatment quarter, averaged across post-treatment quarters.

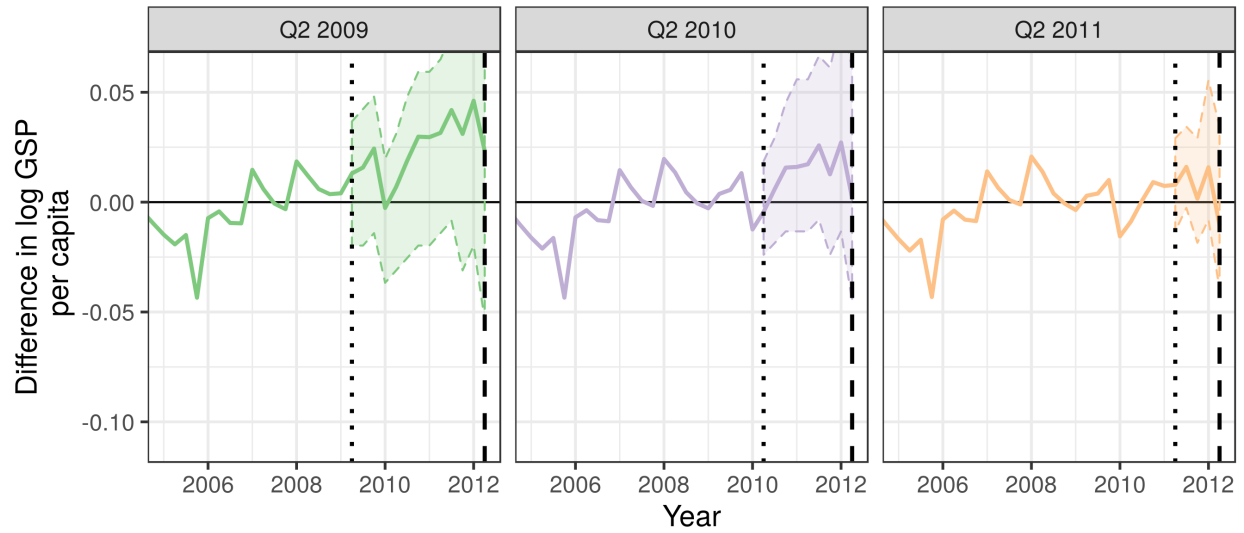


Figure F.6: Placebo point estimates \pm two standard errors for SCM with placebo treatment times in Q2 2009, 2010, and 2011. Scale begins in 2005 to highlight placebo estimates.

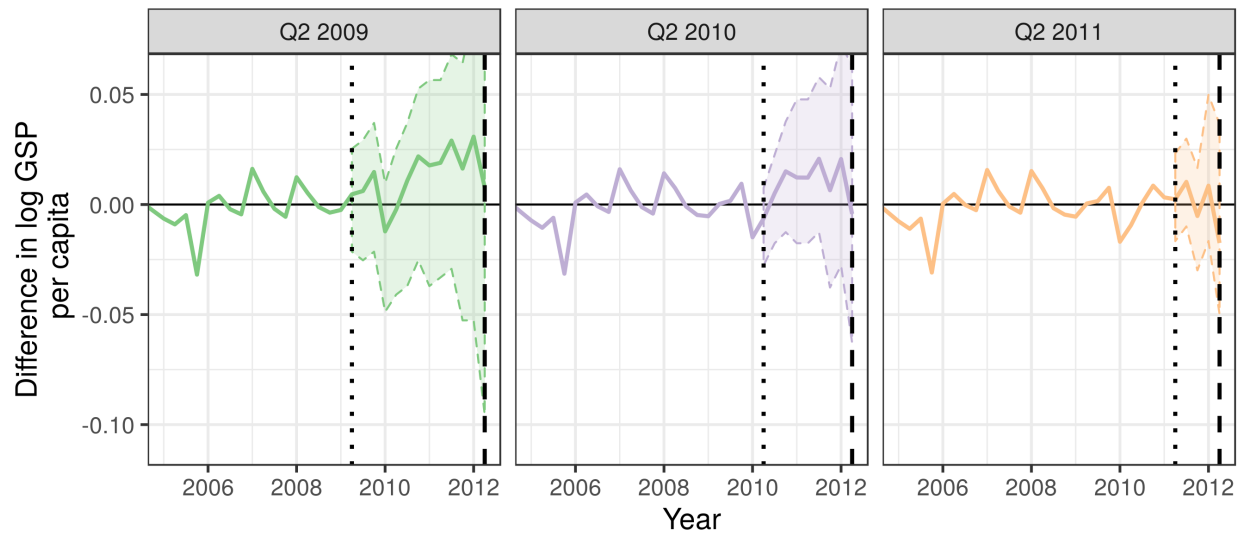


Figure F.7: Placebo point estimates \pm two standard errors for ridge ASCM with placebo treatment times in Q2 2009, 2010, and 2011. Scale begins in 2005 to highlight placebo estimates.

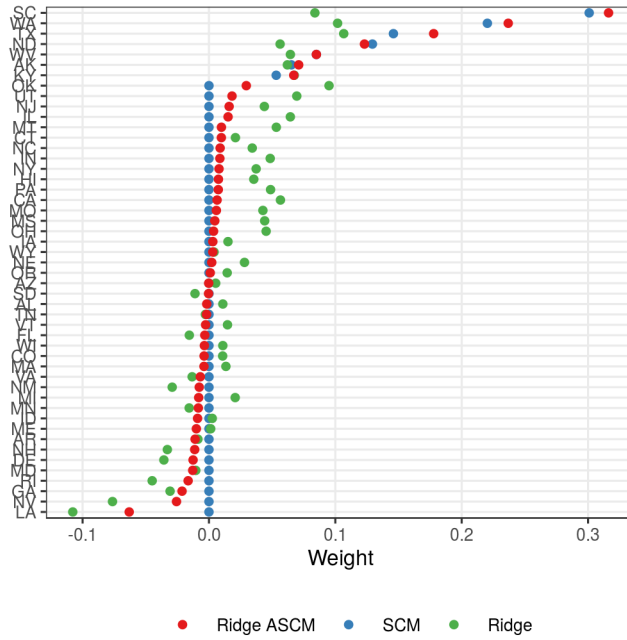


Figure F.8: Donor unit weights for SCM, ridge regression, and ridge ASCM balancing lagged outcomes.

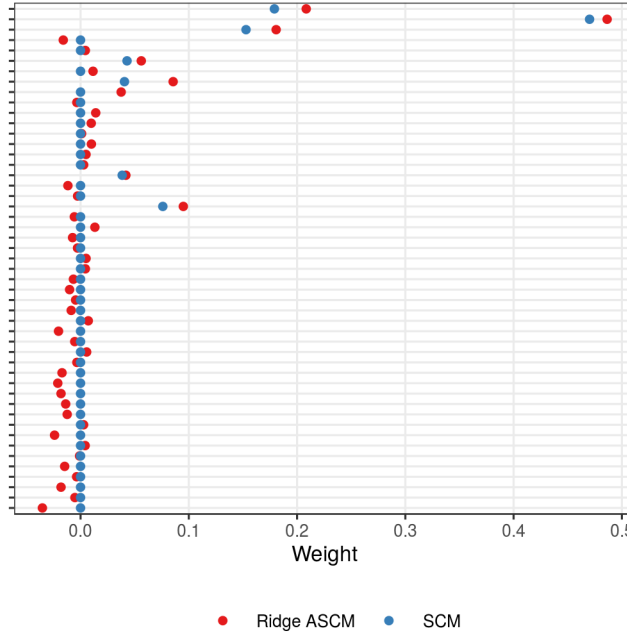


Figure F.9: Donor unit weights for SCM and ridge ASCM fit on lagged outcomes after residualizing out auxiliary covariates.

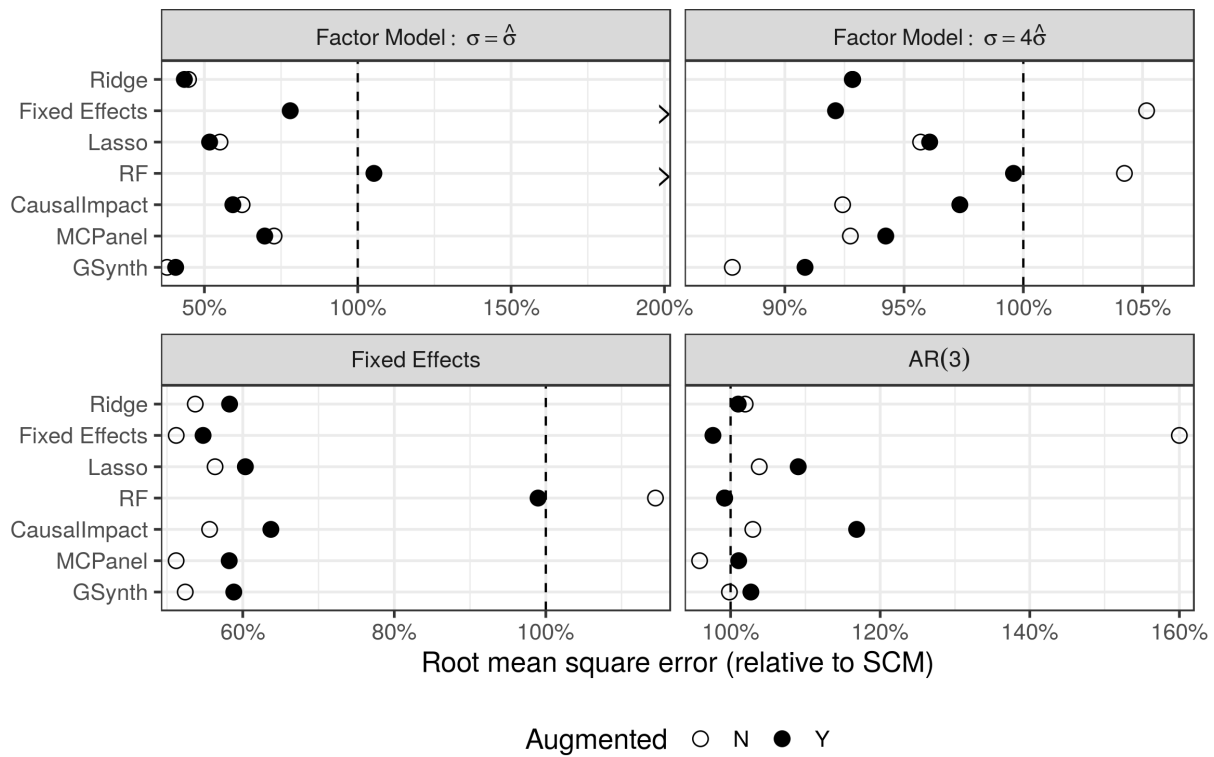


Figure F.10: RMSE for different augmented and non-augmented estimators across outcome models.

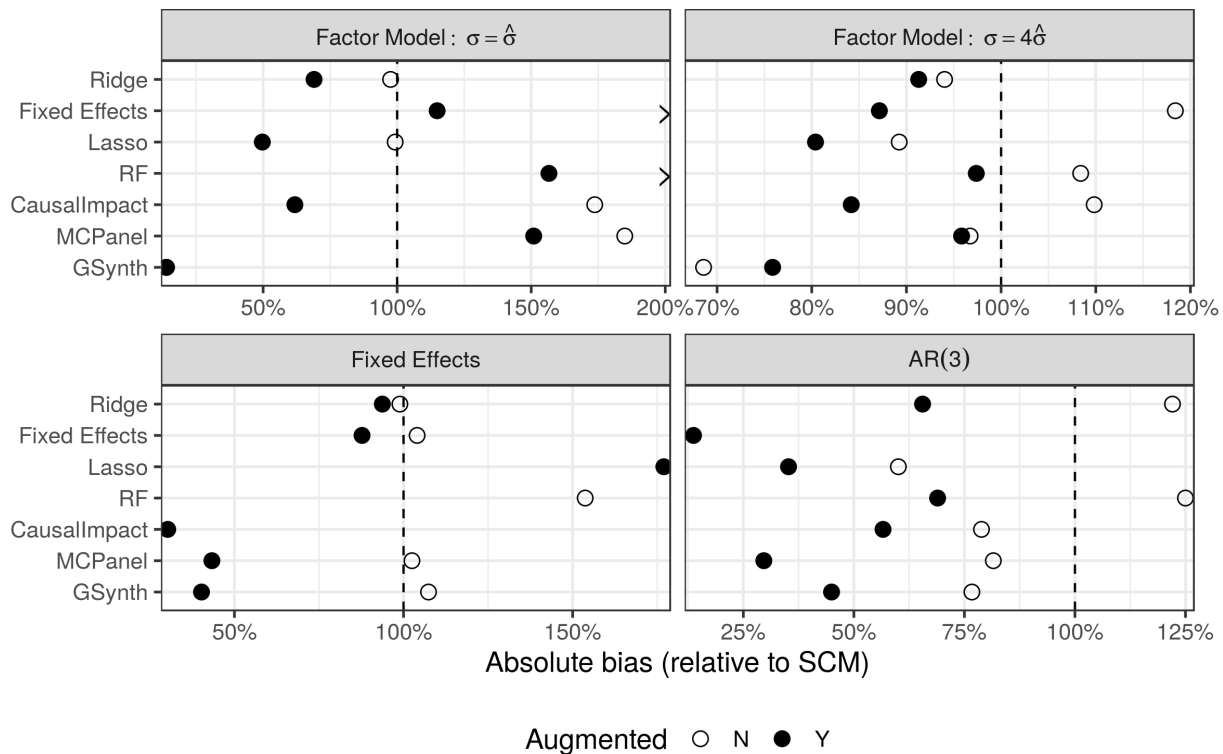


Figure F.11: Bias for different augmented and non-augmented estimators across outcome models conditioned on SCM fit in the top quintile.

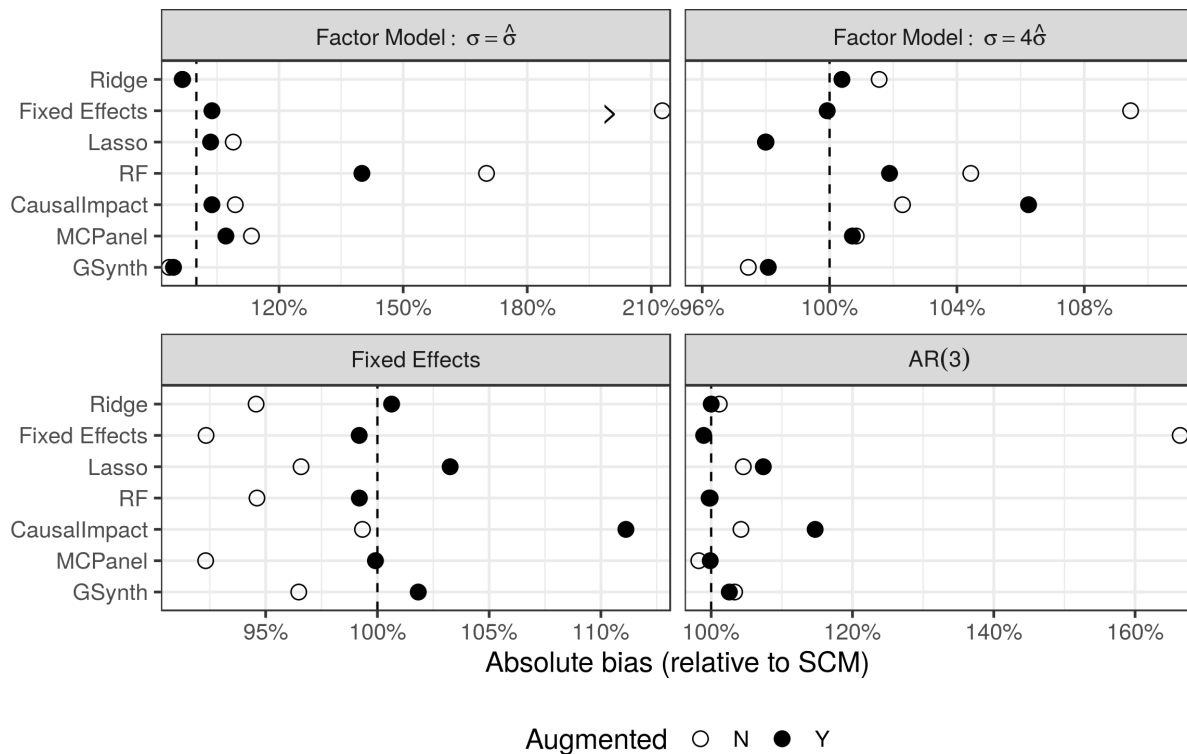


Figure F.12: RMSE for different augmented and non-augmented estimators across outcome models conditioned on SCM fit in the top quintile.

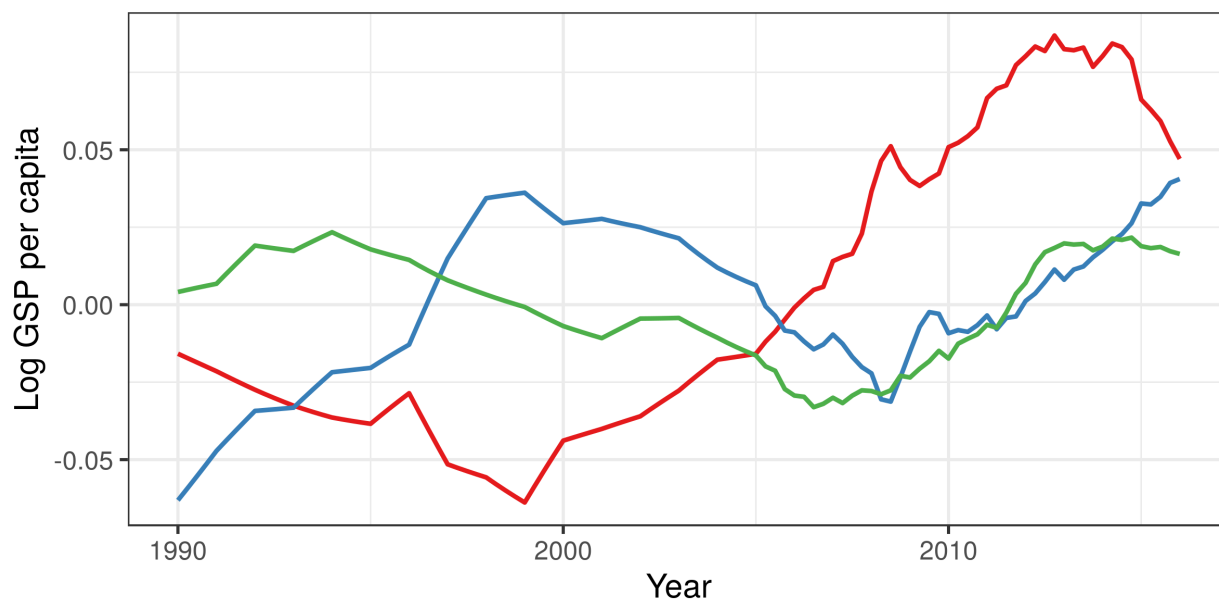


Figure F.13: Latent factors for calibrated simulation studies.

References

- Abadie, A. (2019). Using synthetic controls: Feasibility, data requirements, and methodological aspects. *Journal of Economic Literature*.
- Abadie, A., A. Diamond, and J. Hainmueller (2010). Synthetic Control Methods for Comparative Case Studies: Estimating the Effect of California’s Tobacco Control Program. *Journal of the American Statistical Association* 105(490), 493–505.
- Abadie, A., A. Diamond, and J. Hainmueller (2015). Comparative Politics and the Synthetic Control Method. *American Journal of Political Science* 59(2), 495–510.
- Abadie, A. and J. Gardeazabal (2003). The Economic Costs of Conflict: A Case Study of the Basque Country. *The American Economic Review* 93(1), 113–132.
- Abadie, A. and G. W. Imbens (2011). Bias-corrected matching estimators for average treatment effects. *Journal of Business & Economic Statistics* 29(1), 1–11.
- Abadie, A. and J. L’Hour (2018). A penalized synthetic control estimator for disaggregated data.
- Amjad, M., D. Shah, and D. Shen (2018). Robust synthetic control. *The Journal of Machine Learning Research* 19(1), 802–852.
- Arkhangelsky, D., S. Athey, D. A. Hirshberg, G. W. Imbens, and S. Wager (2019). Synthetic difference in differences. *arXiv preprint arXiv:1812.09970*.
- Athey, S., M. Bayati, N. Doudchenko, G. Imbens, and K. Khosravi (2017). Matrix Completion Methods for Causal Panel Data Models. *arxiv 1710.10251*.
- Athey, S. and G. W. Imbens (2017). The state of applied econometrics: Causality and policy evaluation. *Journal of Economic Perspectives* 31(2), 3–32.
- Athey, S., G. W. Imbens, and S. Wager (2018). Approximate residual balancing: debiased inference of average treatment effects in high dimensions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 80(4), 597–623.
- Bai, J. (2009). Panel data models with interactive fixed effects. *Econometrica* 77(4), 1229–1279.
- Barber, R. F., E. J. Candès, A. Ramdas, and R. J. Tibshirani (2019). Predictive inference with the jackknife+. *arXiv preprint arXiv:1905.02928*.
- Ben-Michael, E., A. Feller, and J. Rothstein (2019). Synthetic controls and weighted event studies with staggered adoption. *arXiv preprint arXiv:1912.03290*.
- Bilinski, A. and L. Hatfield (2020). Goldilocks and the pre-intervention time series. Technical report.
- Botosaru, I. and B. Ferman (2019). On the role of covariates in the synthetic control method. *The Econometrics Journal* 22(2), 117–130.
- Breidt, F. J. and J. D. Opsomer (2017). Model-Assisted Survey Estimation with Modern Prediction Techniques. *Statistical Science* 32(2), 190–205.

- Brodersen, K. H., F. Gallusser, J. Koehler, N. Remy, and S. L. Scott (2015). Inferring Causal Impact using Bayesian Structural Time-Series Models. *The Annals of Applied Statistics* 9(1), 247–274.
- Cassel, C. M., C.-E. Sarndal, and J. H. Wretman (1976). Some results on generalized difference estimation and generalized regression estimation for finite populations. *Biometrika* 63(3), 615–620.
- Cattaneo, M. D., Y. Feng, and R. Titiunik (2019). Prediction intervals for synthetic control methods. *arXiv preprint arXiv:1912.07120*.
- Chattopadhyay, A., Christopher H. Hase, and J. R. Zubizarreta (2020). Balancing Versus Modeling Approaches to Weighting in Practice. *Statistics in Medicine in press*.
- Chernozhukov, V., K. Wuthrich, and Y. Zhu (2018). Inference on average treatment effects in aggregate panel data settings. *arXiv preprint arXiv:1812.10820*.
- Chernozhukov, V., K. Wüthrich, and Y. Zhu (2019). An Exact and Robust Conformal Inference Method for Counterfactual and Synthetic Controls. Technical report.
- Donohue, J. J., A. Aneja, and K. D. Weber (2017). Right-to-carry laws and violent crime: A comprehensive assessment using panel data and a state-level synthetic control analysis. Technical report, National Bureau of Economic Research.
- Doudchenko, N. and G. W. Imbens (2017). Difference-In-Differences and Synthetic Control Methods: A Synthesis. *arxiv 1610.07748*.
- Dube, A. and B. Zipperer (2015). Pooling multiple case studies using synthetic controls: An application to minimum wage policies.
- Ferman, B. (2019). On the Properties of the Synthetic Control Estimator with Many Periods and Many Controls.
- Ferman, B. and C. Pinto (2018). Synthetic controls with imperfect pre-treatment fit.
- Gobillon, L. and T. Magnac (2016). Regional policy evaluation: Interactive fixed effects and synthetic controls. *Review of Economics and Statistics* 98(3), 535–551.
- Hastie, T., J. Friedman, and R. Tibshirani (2009). *The elements of statistical learning*. Springer series in statistics New York.
- Hazlett, C. and Y. Xu (2018). Trajectory balancing: A general reweighting approach to causal inference with time-series cross-sectional data.
- Hirshberg, D. A., A. Maleki, and J. Zubizarreta (2019). Minimax linear estimation of the retargeted mean. *arXiv preprint arXiv:1901.10296*.
- Hirshberg, D. A. and S. Wager (2018). Augmented Minimax Linear Estimation.
- Hsiao, C., Q. Zhou, et al. (2018). Panel parametric, semi-parametric and nonparametric construction of counterfactuals-california tobacco control revisited. Technical report.

- Kellogg, M., M. Mogstad, G. Pouliot, and A. Torgovitsky (2020). Combining matching and synthetic controls to trade off biases from extrapolation and interpolation. Technical report, National Bureau of Economic Research.
- King, G. and L. Zeng (2006). The dangers of extreme counterfactuals. *Political Analysis* 14(2), 131–159.
- Kline, P. (2011). Oaxaca-Blinder as a reweighting estimator. In *American Economic Review*, Volume 101, pp. 532–537.
- Kreif, N., R. Grieve, D. Hangartner, A. J. Turner, S. Nikolova, and M. Sutton (2016). Examination of the synthetic control method for evaluating health policies with multiple treated units. *Health economics* 25(12), 1514–1528.
- Li, K. T. (2017). Statistical Inference for Average Treatment Effects Estimated by Synthetic Control Methods.
- Minard, S. and G. R. Waddell (2018). Dispersion-weighted synthetic controls.
- Neyman, J. (1990 [1923]). On the application of probability theory to agricultural experiments. essay on principles. section 9. *Statistical Science* 5(4), 465–472.
- Powell, D. (2018). Imperfect synthetic controls: Did the massachusetts health care reform save lives?
- Rickman, D. S. and H. Wang (2018). Two tales of two us states: Regional fiscal austerity and economic performance. *Regional Science and Urban Economics* 68, 46–55.
- Robbins, M., J. Saunders, and B. Kilmer (2017). A Framework for Synthetic Control Methods With High-Dimensional, Micro-Level Data: Evaluating a Neighborhood-Specific Crime Intervention. *Journal of the American Statistical Association* 112(517), 109–126.
- Robins, J. M., A. Rotnitzky, and L. P. Zhao (1994). Estimation of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association* 89(427), 846–866.
- Rubin, D. B. (1973). The use of matched sampling and regression adjustment to remove bias in observational studies. *Biometrics*, 185–203.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology* 66(5), 688.
- Rubin, D. B. (1980). Comment on “randomization analysis of experimental data: The fisher randomization test”. *Journal of the American Statistical Association* 75(371), 591–593.
- Samartsidis, P., S. R. Seaman, A. M. Presanis, M. Hickman, D. De Angelis, et al. (2019). Assessing the causal effect of binary interventions from observational panel data with few treated units. *Statistical Science* 34(3), 486–503.
- Soriano, D., E. Ben-Michael, P. Bickel, A. Feller, and S. Pimentel (2020). Sensitivity analysis for balancing weights. Technical report. working paper.

- Tan, Z. (2017). Regularized calibrated estimation of propensity scores with model misspecification and high-dimensional data.
- Toulis, P. and A. Shaikh (2018). Randomization tests in observational studies with time-varying adoption of treatment.
- Vovk, V., A. Gammerman, and G. Shafer (2005). *Algorithmic learning in a random world*. Springer.
- Wainwright, M. (2018). *High dimensional statistics: a non-asymptomatic viewpoint*.
- Wang, Y. and J. R. Zubizarreta (2018). Minimal Approximately Balancing Weights: Asymptotic Properties and Practical Considerations.
- Xu, Y. (2017). Generalized Synthetic Control Method: Causal Inference with Interactive Fixed Effects Models. *Political Analysis* 25, 57–76.
- Zhao, Q. (2018). Covariate Balancing Propensity Score by Tailored Loss Functions. *Annals of Statistics*, forthcoming.
- Zhao, Q. and D. Percival (2017). Entropy balancing is doubly robust. *Journal of Causal Inference* 5(1).
- Zubizarreta, J. R. (2015). Stable Weights that Balance Covariates for Estimation With Incomplete Outcome Data. *Journal of the American Statistical Association* 110(511), 910–922.