



Computing Synthetic Controls Using Bilevel Optimization

Pekka Malo¹ · Juha Eskelinen¹ · Xun Zhou² · Timo Kuosmanen³

Accepted: 6 September 2023
© The Author(s) 2023

Abstract

The synthetic control method (SCM) represents a notable innovation in estimating the causal effects of policy interventions and programs in a comparative case study setting. In this paper, we demonstrate that the data-driven approach to SCM requires solving a bilevel optimization problem. We show how the original SCM problem can be solved to the global optimum through the introduction of an iterative algorithm rooted in Tykhonov regularization or Karush–Kuhn–Tucker approximations.

Keywords Causal effects · Comparative case studies · Policy impact assessment · Bilevel optimization

JEL Classification C31 · C54 · C61

1 Introduction

The synthetic control method (SCM), originally introduced by Abadie and Gardeazabal (2003), is an appealing tool for evaluating the causal treatment effects of policy interventions and programs in comparative case studies (Athey & Imbens, 2017). SCM has been employed in a large number of important applications (e.g., Abadie et al., 2010, 2015; Acemoglu et al., 2016; Cavallo et al., 2013;

✉ Xun Zhou
xun.zhou@york.ac.uk

Pekka Malo
pekka.malo@aalto.fi

Juha Eskelinen
juha.p.eskelinen@aalto.fi

Timo Kuosmanen
timo.kuosmanen@utu.fi

¹ Department of Information and Service Management, Aalto University School of Business, Espoo, Finland

² Department of Environment and Geography, University of York, York, UK

³ Department of Economics, Turku School of Economics, University of Turku, Turku, Finland

Gobillon & Magnac, 2016; Kleven et al., 2013; Bayer & Aklın, 2020). Since the outbreak of the COVID-19 pandemic, SCM has extensively been applied to identify the impacts of the pandemic-related restrictions (e.g., Alfano et al., 2021; Bonander et al., 2021; Cole et al., 2020; Lang et al., 2022; Mills & Rüttenauer, 2022; Mitze et al., 2020; Sehgal, 2021; Xin et al., 2021).

SCM estimates the causal treatment effect by constructing a counterfactual of the treated unit (i.e., synthetic control) using a convex combination of similar units not exposed to the treatment (i.e., donors). The convex combination requires non-negative weights that sum to one to avoid extrapolation. The weights are determined to ensure that the treated unit and the synthetic control resemble each other as closely as possible prior to the treatment, both with respect to the outcome of interest and some observed predictors. Since there are typically multiple predictors, the predictors are also weighted using another set of non-negative weights. Abadie and Gardeazabal (2003) and Abadie et al. (2010) discuss several alternative approaches to specify the predictor weights, including the use of subjective weights. In practice, a majority of published SCM applications resort to a data-driven procedure where the weights of predictors and donors are jointly optimized to minimize the mean squared prediction error of the synthetic control over the pre-treatment period, applying the *Synth* package described in Abadie et al. (2011), which is available for R, Matlab, and Stata.

Despite the popularity of SCM, rather surprisingly, no explicit mathematical formulation of how the predictor weights and the donor weights are jointly optimized has been presented in the literature. Several recent studies note that the synthetic controls produced by standard computational packages available for SCM may encounter numerical instability or fail to achieve the optimum (e.g., Albalade et al., 2021; Becker & Klößner, 2017, 2018; Becker et al., 2018; Klößner, 2015; Kuosmanen et al., 2021).

The purpose of the present paper is to provide a comprehensive investigation into the optimization problem that needs to be solved to compute the synthetic control weights. Unfortunately, the explicit formulation of the SCM problem reveals that computing the synthetic controls entails solving an NP-hard problem, referred to as a bilevel optimization problem (e.g., Hansen et al., 1992; Vicente et al., 1994). In essence, the task of computing synthetic controls turns out to be more challenging than any previous SCM studies recognize. This insight sheds light on the numerical instability reported by Klößner et al. (2015), among others. To address this problem, we develop an iterative algorithm for solving the original SCM problem, based on Tykhonov regularization or Karush–Kuhn–Tucker approximations. We formally prove that the proposed algorithm is guaranteed to converge to the optimal solution.

The rest of the paper is organized as follows. Section 2 introduces the SCM method and formulates the data-driven approach to compute the predictor and donor weights as a bilevel optimization problem. Section 3 develops an iterative algorithm that is guaranteed to converge to the optimal solution. Section 4 applies the proposed algorithm to the data of the seminal SCM application to the California tobacco control program and compares the empirical results with those produced by the existing computational tools for SCM. Section 5 presents our concluding remarks and

discusses potential avenues for future research. Proofs of theorems and the implementation of the descent algorithm are presented in the Appendices.

2 Synthetic Control Method

2.1 Preliminaries

Following the usual notation (e.g., Abadie, 2021), suppose we observe units $j = 1, \dots, J + 1$, where the first unit is exposed to the intervention and the J remaining units are control units that can contribute to the synthetic control. The set of J control units is referred to as the donor pool. For the sake of clarity, we denote the number of time periods prior to treatment as T^{pre} and the number of time periods after the treatments as T^{post} . The outcome of interest is denoted by Y : column vectors Y_1^{pre} and Y_1^{post} with T^{pre} and T^{post} rows, respectively, refer to the time series of the pre-treatment and post-treatment outcomes of the treated unit. Similarly, matrices Y_0^{pre} and Y_0^{post} with J columns refer to the pre-treatment and post-treatment outcomes of the control group, respectively.

Ideally, the impact of treatment could be measured as

$$\alpha = Y_1^{\text{post}} - Y_1^{\text{post,N}}, \tag{1}$$

where $Y_1^{\text{post,N}}$ refers to the counterfactual outcome that would occur if the unit was not exposed to the treatment. If one could observe the outcomes $Y_1^{\text{post,N}}$ in an alternative state of nature, where the unit was not exposed to the treatment, then one could simply calculate the elements of vector α . The main challenge in the estimation of the treatment effect is that only Y_1^{post} is observable, whereas the counterfactual $Y_1^{\text{post,N}}$ is not.

The goal of SCM is to construct a synthetic control group to estimate the counterfactual $Y_1^{\text{post,N}}$. The key idea of the SCM is to use the convex combination of the observed outcomes of the control units Y_0^{post} as an estimator of $Y_1^{\text{post,N}}$. Formally, the SCM estimator is defined as

$$\hat{\alpha} = Y_1^{\text{post}} - Y_0^{\text{post}} W, \tag{2}$$

where the elements of vector W are non-negative and sum to one. The weights W characterize the synthetic control, that is, a counterfactual path of outcomes for the treated unit in the absence of treatment.

To set the weights W , the simplest approach considered by Abadie and Gardeazabal (2003) is to track the observed path of pre-treatment outcomes as closely as possible to minimize the mean squared prediction error (MSPE). That is, one could apply the weights W that solve the following constrained least squares problem

$$\min_{W \in \mathcal{W}} L(W) = \frac{1}{T^{\text{pre}}} \left\| Y_1^{\text{pre}} - Y_0^{\text{pre}} W \right\|^2, \tag{3}$$

where

$$\mathcal{W} = \left\{ W \in \mathbb{R}^J : \sum_{j=2}^{J+1} W_j = 1, W_j \geq 0, j = 2, \dots, J+1 \right\} \quad (4)$$

is the set of admissible weights for control units and $\|\cdot\|$ denotes the usual Euclidean norm. The constraints on the weights W ensure that the synthetic control is a convex combination of the control units in the pool of donors. The fact that SCM does not involve extrapolation is considered one of its greatest advantages over regression analysis (e.g., Abadie, 2021). Note that if we relax the constraints on weights W , then the unconstrained minimization problem reduces to the classic OLS problem without the intercept term. In that case, one could simply regress the time series Y_1^{pre} on the parallel outcomes of the J donors in the control group and set the weights W equal to the corresponding OLS coefficients. While the OLS problem has the well-known closed-form solution that satisfies the first-order conditions, the optimal solution to the constrained least squares problem stated above is typically a corner solution where at least some of the constraints on weights W are binding. The constrained least squares problem can be efficiently solved by quadratic programming (QP) algorithms such as CPLEX, which are guaranteed to converge to the global optimum.

In addition to the outcome of interest, an integral part of SCM is to utilize additional information observed during the pre-treatment period. Suppose we observe K variables referred to as predictors (also known as growth factors, characteristics, or covariates), which are observed prior to the treatment or are unaffected by the treatment, which can influence the evolution of Y . These predictors are denoted by a $(K \times 1)$ vector X_1 and a $(K \times J)$ matrix X_0 , respectively.¹ Abadie et al. (2010) prove unbiasedness and consistency of the SCM under the condition that the synthetic control yields perfect fit to the predictors, that is, $X_1 = X_0W$. Abadie (2021) notes that “*In practice, the condition $X_1 = X_0W$ is replaced by the approximate version $X_1 \approx X_0W$. It is important to notice, however, that for any particular dataset, there are no ex-ante guarantees on the size of the difference $X_1 - X_0W$. When this difference is large, Abadie et al. (2010) recommend against the use of synthetic controls because of the potential for substantial biases.*”

Since the K predictors included in X do not necessarily have the same effect on the outcomes Y , Abadie and Gardeazabal (2003) introduce a $(K \times K)$ diagonal matrix V where the diagonal elements are weights of the predictors that reflect the

¹ A common practice in SCM is to include some convex combinations of the pre-treatment outcomes also as predictors (see Abadie et al., 2010, 2015, for discussion). However, Kaul et al. (2022) demonstrate that including all pre-treatment outcomes as predictors is not a good idea because the predictors become completely redundant in that case.

relative importance of the predictors. The diagonal elements of V must be non-negative² and are usually normalized to sum to unity.³ That is

$$V \in \left\{ \text{diag}(V) : V \in \mathbb{R}^{K \times K}, \sum_{k=1}^K V_{kk} = 1, V_{kk} \geq 0 \right\} =: \mathcal{V}, \tag{5}$$

which is a subset of all non-negative diagonal matrices.

Both Abadie and Gardeazabal (2003) and Abadie et al. (2010) suggest that weights V could be subjectively determined. However, virtually all known applications of SCM resort to the data-driven procedure suggested by the authors. Unfortunately, these seminal papers do not explicitly state the required optimization problem. A closer examination of the SCM problem in the next section reveals that the SCM problem is far from trivial from the computational point of view.

2.2 Bilevel Optimization Problem

Since Abadie and Gardeazabal (2003) and Abadie et al. (2010) only state the SCM problem implicitly, to gain a better understanding of the data-driven approach, the first step is to formulate the SCM problem explicitly. By comparing with the original SCM articles, it is easy to verify that the optimal weights V^*, W^* must be obtained as an optimal solution to the following optimistic bilevel optimization problem (cf. Albalade et al., 2021)

$$\min_{V, W} L_V(V, W) = \frac{1}{T_{\text{pre}}} \|Y_1^{\text{pre}} - Y_0^{\text{pre}} W(V)\|^2 \tag{6}$$

$$\text{s.t. } W(V) \in \Psi(V) := \underset{W \in \mathcal{W}}{\text{argmin}} L_W(V, W) = \|X_1 - X_0 W\|_V^2, \tag{7}$$

$$V \in \mathcal{V},$$

where $\|\cdot\|_V$ is a semi-norm parametrized by V , and $\Psi : \mathcal{V} \rightrightarrows \mathcal{W}$ denotes the solution set mapping from upper-level decisions to the set of global optimal solutions of the lower-level problem. For any $(K \times 1)$ real vector Z , we define $\|Z\|_V = (Z^T V Z)^{1/2}$.

This becomes a proper norm only when V is positive-definite. If we denote the diagonal elements of V by v_1, \dots, v_K , we can write the lower-level objective as

$$L_W(V, W) = \sum_{k=1}^K v_k \left(X_{k,1} - \sum_{j=2}^{J+1} X_{k,j} W_j \right)^2,$$

² While Abadie et al. (2010) assume that the diagonal elements must be positive, a positive real number can be arbitrarily close to zero, and therefore, the distinction between positive and non-negative model variables has no real meaning in optimization unless one imposes some explicit lower bound, e.g., $V_{kk} \geq 0.01$. Becker and Klößner (2018) set a lower bound $V_{kk} \geq 0.00000001$, which is so low that it has no practical meaning.

³ Of course, other normalizations are possible, but we here restrict attention to the most standard normalization that allows one to interpret the diagonal elements of V as shared weights that sum to one.

which allows the lower-level problem to be interpreted as an importance-weighted least squares with weight constraints. As pointed out by Klößner and Pfeifer (2015), this original setup can be easily extended to allow the treatment of predictor data as time series, while maintaining the original structure of the optimization problem.

The explicit formulation of the optimization problem reveals several points worth noting. First, the SCM problem is a bilevel optimization problem, which is far from trivial from the computational point of view. The minimization problem (7) referred to the lower-level problem, and problem (6) is called the upper-level problem; the SCM literature commonly uses the terms inner and outer problems, but the meaning is the same. The problem is solvable when it is interpreted as an optimistic bilevel problem, but the global optimum is not necessarily unique.

Proposition 1 *The synthetic control problem defined by (6)–(7) has a global optimistic solution $(\bar{V}, \bar{W}) \in \mathcal{V} \times \mathcal{W}$.*

Unfortunately, the bilevel optimization problems are generally NP-hard (Hansen et al., 1992; Vicente et al., 1994). In particular, the hierarchical optimization structure can introduce difficulties such as non-convexity and disconnectness (e.g., Sinha et al., 2013), which are also problematic in the present setting, as will be demonstrated in the next section.

Second, the explicit statement of the optimization problem makes it more evident that the optimal solution will typically be a corner solution where at least some of the first-order conditions do not hold. This causes a serious problem for the usual derivative-based optimization tools. This observation can help to explain at least partly the numerical instability of the SCM results, observed by Becker and Klößner (2017) and Klößner et al. (2015), among others. The general-purpose computational tools are simply ill-equipped for the task at hand. If the weights W, V are arbitrarily determined by an ad hoc computational tool that fails to converge to a feasible and unique global optimum, then all attractive theoretical properties of the estimator are no longer guaranteed.

3 Iterative Algorithm

The purpose of this section is to discuss a general algorithm for solving the original SCM problem (4)–(5) where the predictor weights V are jointly optimized with the donor weights W . Since the general algorithm proves computationally demanding, we start by checking whether the unconstrained SCM problem (3) is a feasible solution as well as the possibility of corner solutions. It is noteworthy that surprisingly many of the SCM problems encountered in practice admit either an unconstrained solution or a corner solution. In case the optimal solution is not found through these feasibility checks, we suggest continuing the search for an optimal solution using a descent algorithm based on the Tykhonov regularization technique or Karush–Kuhn–Tucker (KKT) approximations.

To highlight the importance of coordination between the upper-level and lower-level problems, we can rephrase the lower-level problem (7) as

$$\min_{W \in \mathcal{W}} L_W^\varepsilon(V, W) = \frac{1}{K} \|X_1 - X_0 W\|_{V^*}^2 + \varepsilon \|Y_1^{\text{pre}} - Y_0^{\text{pre}} W(V)\|^2 \tag{8}$$

where $\varepsilon > 0$ denotes an infinitesimally small non-Archimedean scalar.⁴ Introducing the upper-level objective as a part of the lower-level QP problem in (8) makes a subtle but important difference compared to problem (7): the primary objective of both (7) and (8) is to minimize the loss function L_W with respect to predictors X . However, if there are alternate optima W^* that minimize the loss function L_W , problem (8) will choose the best solution for the upper-level problem.

Proposition 2 *For a given set of weights V^* , let $W_\varepsilon(V^*)$ denote the unique optimal solution to problem (8) for any $\varepsilon > 0$. Then, we have that*

$$\lim_{\varepsilon \rightarrow 0^+} W_\varepsilon(V^*) \in \underset{W}{\operatorname{argmin}}\{L_V(V^*, W) : W \in \Psi(V^*)\}.$$

The proof of the proposition is simple and can be omitted. Having ensured that constraint (5) holds, it is important to note that the optimal weights W that minimize $\|X_1 - X_0 W\|_{V^*}^2$ need not be unique. This is particularly relevant when there exist W that satisfy $\|X_1 - X_0 W\|_{V^*}^2 = 0$. In such cases, the non-Archimedean ε plays an important role by allowing us to select among the alternate optima for (5) the optimal weights W to minimize the upper-level objective (6).

Proposition 2 provides a useful result for SCM applications where the weights V are given. Recall that weights V might be subjectively determined, as Abadie and Gardeazabal (2003) and Abadie et al. (2010) suggest. Proposition 2 also demonstrates the critical importance of introducing an explicit link between the lower-level problem and the upper-level problem. In general, there can be many alternate optima where the loss function goes to zero, $L_W = 0$. Without coordination, there is no guarantee that the SCM package would converge to the optimum. The lack of an explicit link between the upper-level and the lower-level problem is the most fundamental reason why the *Synth* package fails to reach the optimum.

3.1 Checking the Feasibility of an Unconstrained Solution

Consider first the situation where no predictors are used (i.e., $K = 0$). In this case, the bilevel optimization problem (6)–(7) reduces to the constrained regression problem (3). Problem (3) has a quadratic objective function and a set of linear constraints, which guarantees the existence of a unique global optimum when the usual assumptions of regression analysis hold (i.e., no rank deficiency). Such quadratic programming problems are considered straightforward from the computational point of view. While general-purpose derivative-based tools may struggle with the

⁴ The use of non-Archimedean ε was introduced by Charnes (1952) to avoid degeneracy in linear programming.

constraints, the simplex-based algorithms (e.g. the CPLEX solver) will converge to the global optimum.

Let $L(W^{**}) = \min_{W \in \mathcal{W}} L(W)$ denote the optimal solution to the problem (3), which is unique when no rank deficiency is present. As Kaul et al. (2022) correctly note, this solution is the lower bound for the optimal solution to the problem (6):

$$L_V(V, W) \geq L(W^{**}) \text{ for all } V \in \mathcal{V}, W \in \mathcal{W}. \tag{9}$$

Intuitively, imposing additional constraints can never improve the optimal solution. To test if there exist importance weights $V \in \mathcal{V}$ such that W^{**} is a feasible solution to the lower-level problem (7), we next solve the following linear programming (LP) problem

$$\min_{V \in \mathcal{V}} L_W(V, W^{**}) = (X_1 - X_0 W^{**})^\top V (X_1 - X_0 W^{**}). \tag{10}$$

While the objective function of problem (10) is the same as that of the lower-level problem (7) in that both problems minimize the same loss function, problem (7) is minimized with respect to weights W , whereas problem (10) is minimized with respect to weights V , taking W^{**} as given. This LP problem finds the optimal predictor weights V to support the relaxed problem (3). Denote the optimal solution to problem (10) as V^{**} . If $L_W(V^{**}, W^{**}) = 0$, the optimal solution has been found. In other words, there exists matrix $V^{**} \in \mathcal{V}$ such that W^{**} is a feasible solution to the lower-level problem (7), i.e. $W^{**} \in \Psi(V^{**})$. Hence, this is also the optimal solution to the bilevel optimization problem (6)–(7).

3.2 Establishing an Upper Bound for L_V

In the context of SCM, the domain of predictor weights V has K basic solutions, with the following diagonal elements: $V_1 = (1, 0, \dots, 0), V_2 = (0, 1, \dots, 0), \dots, V_K = (0, 0, \dots, 1)$. That is, we assign all weight to just one of the predictors and leave zero weight to all other predictors. We can insert the basic solution $V_k, k = 1, \dots, K$ as the weights V in problem (8), and solve the QP problem to find the optimal W_k for each $k = 1, \dots, K$. For each candidate weights $W_k, k = 1, \dots, K$, we calculate the value of the upper-level loss function L_V stated in (6). Finally, we select the basic solution s in $1, \dots, K$ that minimizes L_V . If $L_W(V_s, W_s) = 0$ and $L_V(V_s, W_s) = L(W^{**})$, then the corner solution (V_s, W_s) is one of the optimal solutions. If only $L_W(V_s, W_s) = 0$ but $L_V(V_s, W_s) > L(W^{**})$, the corner solution can be viewed as an upper bound for the optimal value.

Proposition 3 *If there exist weights $(\tilde{V}, \tilde{W}) \in \mathcal{V} \times \mathcal{W}$ satisfying $X_{0k} \tilde{W} = x_{1k}$ for some predictor k , then there exists another feasible solution (V_k, \tilde{W}) for the SCM problem (6)–(7), where $V_k \in \mathcal{V}$ is a corner solution satisfying $L_W(V_k, \tilde{W}) = 0$. If (\tilde{V}, \tilde{W}) is an optimal solution, then also (V_k, \tilde{W}) is an alternative optimal solution for the SCM problem.*

This result demonstrates that whenever the donor weights W satisfy the basic condition required for the consistency of the SCM, $X_1 = X_0W$, even just for a single predictor k , then it is easy to generate feasible solution candidates that are obtained by considering corner solutions with respect to predictor weights V . Intuitively, when the number of predictors is large, it is practically impossible to construct a convex combination of control units that matches the treated unit; in other words, no matrix W that satisfies $X_0W = X_1$ exists. But if we use weights V to reduce the dimensionality of X by assigning some of the predictors a zero weight, then it becomes considerably easier to find vectors W that satisfy $x_{0k}W = x_{1k}$ at least for some predictor k (note x_{0k} is the k th row of matrix X_0 and x_{1k} is a scalar). Consequently, the set of feasible solutions for the SCM problem often contains several candidate solutions that “switch off” the constraints concerning predictors X by assigning zero weight, except for a single predictor k for which a perfect fit is possible. Therefore, it is understandable that many ad hoc tools attempting to solve the SCM problem (6)–(7) may end up assigning all weight to the most favorable predictor and discard all other predictors by assigning the zero weight. These observations can help to explain the empirical observation that the predictors often turn out to have little impact on the synthetic control, which has been noted by several authors (e.g., Ben-Michael et al., 2021; Doudchenko & Imbens 2017; Kaul et al., 2022). While these solutions may not necessarily be optimal for the SCM problem, they can still provide good approximations for the optimal value of the upper-level objective. Note that the previous iterations provide us with the corner solution (V_k, W_k) and the unconstrained solution W^{**} , which can be used for constructing the following bounds for the loss function of the true optimum (V^*, W^*) :

$$L_V(V_s, W_s) \geq L_V(V^*, W^*) \geq L(W^{**}).$$

If the margin of L_V is small and $W_s \approx W^{**}$ by reasonable tolerance, there is no need to iterate further. But if there is a significant gap, the following iterative procedure is guaranteed to find the optimum.

3.3 Finding an Optimal Solution Using Tykhonov Regularization

Building on Proposition 2, the basic idea is to construct an iterative descent algorithm to find the bilevel optimal solution by using the following regularized lower-level problem:

$$\min_{W \in \mathcal{W}} L_W^\epsilon(V, W) = L_W(V, W) + \epsilon L_V(V, W), \tag{11}$$

where $\epsilon > 0$. Note that problem (11) is just a re-stated version of the QP problem (8) above. When the optimal solution to the upper-level problem is uniquely defined, the regularized lower-level problem has considerably better regularity properties than the original formulation. In the literature on bilevel programming, this approach is

known as Tykhonov regularization (Dempe, 2002). By requiring positive definiteness in the upper-level problem, we can make relatively strong claims regarding the properties of the optimal solutions for the regularized problem. Specifically, it can be shown that the unique optimal solution function to the problem (11), denoted by $W_{\epsilon_k}^*(V)$, is Lipschitz continuous and directionally differentiable.

Definition 1 (Lipschitz continuity) A function $z : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is called locally Lipschitz continuous at a point $x^0 \in \mathbb{R}^n$ if there exists an open neighborhood $U_\epsilon(x^0)$ of x^0 and a constant $l < \infty$ such that

$$\|z(x) - z(x')\| \leq l\|x - x'\| \quad \forall x, x' \in U_\epsilon(x^0).$$

Definition 2 (Directional differentiability) A function $z : \mathbb{R}^n \rightarrow \mathbb{R}$ is directionally differentiable at x^0 if for each direction $r \in \mathbb{R}^n$ the following one-sided limit exists:

$$z'(x^0; r) = \lim_{t \rightarrow 0^+} t^{-1}[z(x^0 + tr) - z(x^0)].$$

The value $z'(x^0; r)$ is called the directional derivative of z at $x = x^0$ in direction r .

Proposition 4 Consider the synthetic control problem in (6)–(7) and let the upper-level cross-product matrix $Y_0^T Y_0$ be positive definite. Take any sequence of positive numbers $\{\epsilon^k\}_{k=1}^\infty$ converging to $0+$. Then,

1. the optimal value of the regularized bilevel problem converges to the optimal value of the original problem as $k \rightarrow \infty$ i.e.

$$\min_{V, W} \{L_V(V, W) : W \in \Psi_{\epsilon_k}(V), V \in \mathcal{V}\} \rightarrow L_V^*,$$

where

$$\Psi_{\epsilon_k}(V) = \operatorname{argmin}_{W \in \mathcal{W}} L_W^\epsilon(V, W),$$

$$L_V^* = \min_{V, W} \{L_V(V, W) : W \in \Psi(V), V \in \mathcal{V}\}$$

denote the optimal solution set mapping for (11) and the upper-level optimal value of the original problem, respectively.

2. for each ϵ_k , the unique optimal solution to the regularized lower-level problem (11), denoted by $W_{\epsilon_k}^*(V) \in \Psi_{\epsilon_k}(V)$, is directionally differentiable and

$$\lim_{k \rightarrow \infty} \{W_{\epsilon_k}^*(V)\} = \operatorname{argmin}_W \{L_V(V, W) : W \in \Psi(V)\}$$

for every fixed $V \in \mathcal{V}$.

Based on this result, solving the synthetic control problem is equivalent to considering a sequence of problems

$$\min_V \{L_{\varepsilon_k}(V) : V \in \mathcal{V}\} \text{ for } \varepsilon_k \rightarrow 0+, \quad (12)$$

where the implicitly defined objective function $L_{\varepsilon_k}(V) = L_V(V, W_{\varepsilon_k}^*(V))$ is directionally differentiable with respect to V . The implementation of the descent algorithm is discussed in Appendix B.1. As an alternative to the Tykhonov algorithm, the problem can be also solved using a recently developed approach based on KKT conditions for bilevel problems (Dempe & Franke, 2019). This alternative is briefly described in Appendix B.2.

To summarize this section, the good news is that the SCM problem (6)–(7) is solvable. The bad news is that the required computations prove much more demanding than the original SCM studies assumed. Worse yet, the optimal solution is often a corner solution where most predictors are assigned a zero weight or have a negligible impact. We stress that imposing some small bounds for V (e.g., $V_{kk} \geq 0.01$) would have little impact in practice; the corner solution would simply assign the minimum weight to all predictors, except for the most favorable predictor that would get the maximum weight ($= 1 - 0.01(K - 1)$).

4 Empirical Comparisons

Applying the iterative algorithm proposed in Sect. 3 to the data of the seminal SCM application to the California tobacco control program (Abadie et al., 2010),⁵ we empirically verify that the optimal solution in this original case is indeed a corner solution. Table 1 reports the loss function values of the upper-level problem (L_V) and the lower-level problem (L_W) as well as the donor weights (W) and the predictor weights (V) estimated by different SCM packages available for R.

The corner solution is found superior to the solutions obtained by the standard implementation of *Synth* package described in Abadie et al. (2011)⁶ and the MSCMT (Multivariate Synthetic Control Method using Time Series) package proposed by Becker and Klößner (2018). This observation demonstrates that the existing SCM packages fail to find the optimal solution even in one of the original applications of SCM, which is also used as one of the examples to demonstrate the *Synth* package.

Recall that the value of L_V measures how well the synthetic control matches the pre-treatment outcomes of the treated unit, and this is the upper-level objective to be minimized. In this respect, all computational packages come relatively close to the

⁵ The R code to implement the interactive algorithm is available as online supplementary material at <https://github.com/Xun90/SCM-Debug.git>. The original data of the California application are embedded in the Matlab implementation of *Synth* available at <https://web.stanford.edu/~jhain/synthpage.html>.

⁶ In addition to the standard *Synth* command, we have also considered the `genoud()` option available in *Synth*, as noted in Abadie et al. (2011). However, the use of the `genoud()` option does not improve the matter; in fact, the solution is only worse.

Table 1 California tobacco control application revisited: donor weights, predictor weights, loss functions, and empirical fit by different algorithms

	<i>Synth</i>	MSCMT	Optimum
<i>W</i>			
Utah	0.3432	0.3351	0.3939
Nevada	0.2358	0.2356	0.2049
Montana	0.1820	0.2019	0.2318
Colorado	0.1747	0.1595	0.0148
Connecticut	0.0624	0.0679	0.1091
New Hampshire	0.0000	0.0000	0.0454
<i>V</i>			
Income per capita	0.0006	0.0000	0
Retail price of cigarettes	0.0312	0.3333	0
Population aged 15–19 (%)	0.0034	0.3333	0
Beer consumption per capita	0.0124	0.0000	0
Cigarette sales per capita 1988	0.0682	0.0000	0
Cigarette sales per capita 1980	0.3917	0.0000	1
Cigarette sales per capita 1975	0.4925	0.3333	0
L_V	3.20908	3.07666	2.74366
L_W	0.00170	0.00000	0.00000
R^2	0.97518	0.97621	0.97878

global optimum. It is worth noting that the magnitude of L_V is contingent upon the measurement units of outcomes: for example, multiplying Y_1^{pre} and Y_0^{pre} by 1 thousand would increase L_V by a factor of 1 million. Therefore, it is helpful to measure empirical fit with respect to the pre-treatment outcomes in terms of the coefficient of determination (R^2)—after all, the upper-level problem is just constrained least squares regression without intercept. Such a comparison reveals that the differences in empirical fit are rather marginal; the R^2 statistic varies between 0.97518 (*Synth*) and 0.97878 (optimum). In contrast, the differences in weights W and V across different computational packages are rather dramatic. The results of Table 1 help to illustrate that good empirical fit may be achieved with a wide variety of weights W and V , but there is only one unique global optimum.

The value of L_W measures how well the synthetic control matches the predictors X_1 . While the minimization of L_W is the lower-level objective, the consistency of SCM depends on the (nearly) perfect match with the predictors. In this regard, the value of L_W approaches zero at the global optimum, suggesting a perfect match in terms of the weighted predictors. In contrast, the relatively high value of L_W given by the standard *Synth* command points to the fact that *Synth* fails to converge to the global optimum in the California example. Furthermore, the MSCMT procedure greatly improves L_W in this case and converges to the global optimum. However, the optimal solution is a corner solution that assigns all weight to a single predictor: cigarette sales per capita in 1980 in the California tobacco control application (see

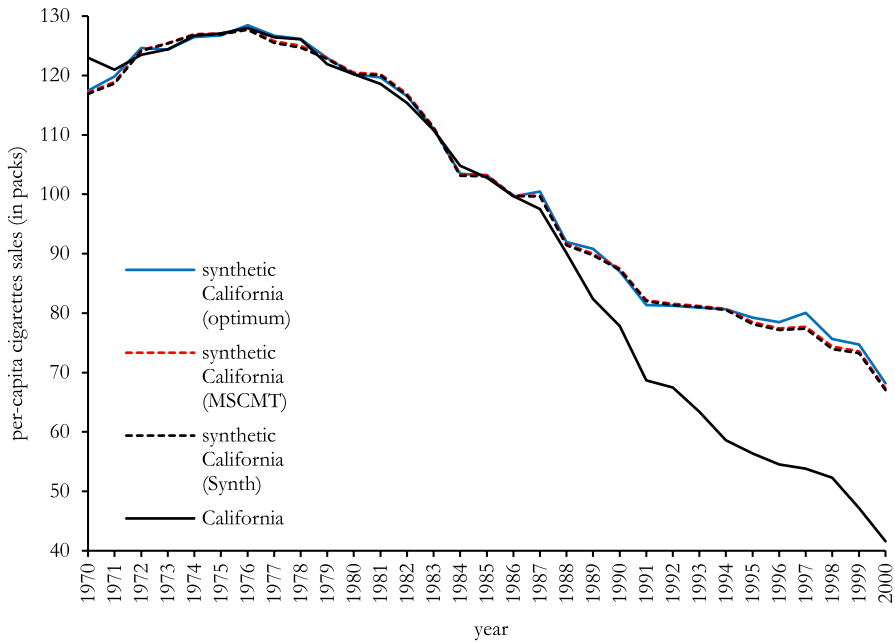


Fig. 1 The impact of suboptimal W weights on the evolution of synthetic California

Table 1). The MSCMT package allocates the weight evenly across three predictors, while the *Synth* package appears to use more balanced weights for predictors; however, note that *Synth* also assigns almost 90% of the predictor weight to cigarette sales per capita (the outcome variable) during two years of the pre-treatment period. Unfortunately, the *Synth* package proves inadequate in solving the optimization problem it is supposed to solve; its predictor weights are not what they are claimed to be, but just artifacts of a computational failure.

Figure 1 illustrates the impact of suboptimal donor weights on the evolution of synthetic California. Fortunately, the qualitative conclusions of this original and highly influential application remain, although the use of the suboptimal weights results in a reduced treatment effect.

5 Conclusions

SCM has proved a highly appealing approach to estimating causal treatment effects within the context of comparative case studies, as demonstrated by numerous published applications. Unfortunately, the standard computational packages aimed at jointly solving the donor weights and the predictor weights have proved numerically unstable. The explicit formulation of the SCM problem as an optimistic bilevel optimization problem highlights that the SCM problem is far from trivial from the computational perspective: the SCM problem is generally NP-hard, significantly exceeding the scope of the computational packages currently in use.

The main contribution of our paper was the introduction of an iterative computational algorithm for solving the original SCM problem. We were the first ones to prove that our SCM algorithm converges to the optimal solution under relatively mild assumptions. This underscores the existence of a theoretically valid approach for solving the SCM problem. However, the optimal solutions to the original SCM formulation are still typically obtained as corner solutions, where most of the predictors carry zero weight. Thus, in practice, it is rarely necessary to apply Tykhonov regularization or KKT approximations to locate the optimal solution. Instead, an optimal solution is usually identified already during the early stages of the iterative procedure.

The computational difficulties of the original SCM formulation do not diminish the conceptual allure of synthetic controls. While we do recognize the value of the data-driven approach to weight determination, it remains crucial to ensure the optimality of the synthetic controls, rather than allowing them to be artifacts of a suboptimal computational tool.

Our findings open various avenues for future research, encompassing both empirical and methodological studies. From the empirical point of view, it would be interesting to apply the proposed algorithm to replicate published SCM studies in order to examine the potential impacts of suboptimal weights on the qualitative conclusions. Becker and Klößner (2017) is an excellent example of such a replication study. We hope that the qualitative results of the influential SCM studies prove robust to the optimization errors that are evidently present, yet this remains to be tested empirically.

From the methodological point of view, the joint optimization of the predictor weights and the donor weights calls for further examination. In particular, the loss function to be minimized requires careful reconsideration to ensure that the optimal solution is reasonable for the intended purposes of using the predictors and that the problem remains computationally tractable. One possibility could involve adopting stepwise optimization of the predictor weights and donor weights, such that the predictor weights are first determined based on alternative criteria (e.g., regression analysis) and subsequently the donor weights are optimized taking the predictor weights as given. We leave this as a fascinating avenue for future research.

Finally, we hope that the insights of our paper could potentially foster further integration of SCM with other estimation approaches such as the difference-in-differences, panel data regression, and machine learning; several recent studies (e.g., Abadie, 2020; Amjad et al., 2018; Arkhangelsky et al., 2021; Ben-Michael et al., 2021; Doudchenko & Imbens, 2017; Xu, 2017) have made impressive progress in this direction.

Appendix A: Proofs of Theorems

A.1 Regularity Conditions for Parametric Optimization

In this section, we will briefly review a few central concepts from parametric optimization literature that we will later need while discussing the notions of optimality for the synthetic control problem. Without loss of generality, the lower-level problem can be stated as a parametric optimization problem

$$\min_y \{f(x, y) : g(x, y) \leq 0, h(x, y) = 0\}, \tag{13}$$

where $f : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$, $g : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^p$, $h : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^q$. The constraints

$$g(x, y) = (g_1(x, y), \dots, g_p(x, y))^T, \\ h(x, y) = (h_1(x, y), \dots, h_q(x, y))^T,$$

are assumed to be smooth vector-valued functions. The problem is a convex parametric optimization problem, when all functions $f(x, \cdot)$, $g_i(x, \cdot)$, $i = 1, \dots, p$, are convex and the functions $h_j(x, \cdot)$, $j = 1, \dots, q$, are affine-linear on \mathbb{R}^n for each fixed $x \in \mathbb{R}^n$. The solution set mapping $\Psi : \mathbb{R}^n \rightrightarrows \mathbb{R}^m$ is defined by

$$\Psi(x) = \underset{y}{\operatorname{argmin}} \{f(x, y) : g(x, y) \leq 0, h(x, y) = 0\},$$

which is a point-to-set mapping from the upper-level decisions to the set of global optimal solutions of the parametric problem. For convex problems, the solution sets $\Psi(x)$ are closed and convex subsets of \mathbb{R}^m .

When it comes to regularity conditions in bilevel programming, the following two conditions have often been utilized. The first condition is concerned with the compactness of the feasible set of the lower-level problem:

Definition 3 (C) The set $\{(x, y) : \mathbb{R}^n \times \mathbb{R}^m : g(x, y) \leq 0, h(x, y) = 0\}$ is non-empty and compact.

This is enough to guarantee that the set of optimal solutions for the parametric problem

$$\Psi(x) := \underset{y}{\operatorname{argmin}} \{f(x, y) : g(x, y) \leq 0, h(x, y) = 0\}$$

is non-empty and compact for each $x \in \{z : \Omega(z) \neq \emptyset\}$, where

$$\Omega(x) = \{y \in \mathbb{R}^m : g(x, y) \leq 0, h(x, y) = 0\}$$

is the feasible set mapping for the lower-level problem.

The second regularity condition is the commonly applied Mangasarian-Fromowitz constraint qualifications:

Definition 4 (MFCQ) We say that Mangasarian-Fromowitz constraint qualification is satisfied at point (x^0, y^0) if there exists a direction $d \in \mathbb{R}^m$ such that

$$\nabla_y g_i(x^0, y^0)d < 0, \text{ for each } i \in I(x^0, y^0) = \{j : g_j(x^0, y^0) = 0\}, \\ \nabla_y h_j(x^0, y^0)d = 0, \text{ for each } j = 1, \dots, q$$

and the gradients of the equality constraints $\{\nabla_y h_j(x^0, y^0) : j = 1, \dots, q\}$ are linearly independent.

These regularity conditions play an important role in ensuring the existence of optimal solutions for optimistic bilevel problems such as the synthetic control problem discussed in this paper. Let $F : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ denote the upper-level objective function that is minimized with respect to upper-level constraints $X := \{x : G(x) \leq 0\}$, $G : \mathbb{R}^n \rightarrow \mathbb{R}^l$. An optimistic solution to a bilevel problem can then be defined as a point solving the following minimization problem:

$$\min_x \{\varphi_0(x) : x \in X\}, \quad (14)$$

where $\varphi_0(x) = \min_y \{F(x, y) : y \in \Psi(x)\}$.

Theorem 1 (Dempe, 2002) *Let the assumptions (C) and (MFCQ) be satisfied at all points $(x, y) \in X \times \mathbb{R}^m$ with $y \in \Omega(x)$. Then, a global solution of the bilevel problem (14) exists provided there is a feasible solution.*

In addition to the existence of optimal solutions, the regularity conditions imply upper-semicontinuity of the optimal solution set mapping.

Definition 5 (Upper semicontinuity) A set-valued mapping $\Psi : \mathbb{R}^n \rightrightarrows \mathbb{R}^m$ is said to be upper semicontinuous at a point $x \in \mathbb{R}^n$ if, for each open set V with $\Psi(x) \subset V$, there exists an open neighborhood $U_\delta(x)$ of x such that $\Psi(x') \subset V$ for each $x' \in U_\delta(x)$.

In the special case, where Ψ is a single-valued mapping, the notion of upper semicontinuity corresponds to the usual continuity of a function.

Theorem 2 (Bank et al., 1982; Dempe, 2002) *Consider the parametric optimization problem (13) at $x = x^0 \in \mathbb{R}^n$ and let the assumptions (C) and (MFCQ) be satisfied for all feasible points (x, y) with $x = x^0$ and $y \in \Omega(x^0)$. Then, the solution set mapping Ψ is upper semicontinuous and the optimal value function φ is continuous at x^0 .*

While the solution set mapping is upper semicontinuous under these relatively weak regularity conditions, it is generally not continuous. The continuity of a solution set mapping is possible only under considerably stronger assumptions such as the strong sufficient optimality condition of second order (SSOC) and constant rank constraint qualification (CRCQ).

Definition 6 (SSOC) The strong sufficient optimality condition of second order holds at (x^0, y^0) if for each pair of Lagrange multipliers $(\lambda, \mu) \in \Lambda(x^0, y^0)$ and for each direction $d \neq 0$ with

$$\begin{aligned} \nabla_y g_i(x^0, y^0)d &= 0, \quad \forall i \in J(\lambda) := \{j : \lambda_j > 0\}, \\ \nabla_y h_j(x^0, y^0)d &= 0, \quad j = 1, \dots, q \end{aligned}$$

we have that

$$d^T \nabla_{yy} L(x^0, y^0, \lambda, \mu) d > 0.$$

Definition 7 (CRCQ) The constant rank constraint qualification holds at point (x^0, y^0) if there exists an open neighborhood $U_\epsilon(x^0, y^0)$ of (x^0, y^0) such that for each subset

$$I \subset I(x^0, y^0) := \{i : g_i(x^0, y^0) = 0\}, J \subset \{1, \dots, q\},$$

the family of gradient vectors

$$\{\nabla_y g_i(x, y) : i \in I\} \cup \{\nabla_y h_j(x, y) : j \in J\}$$

has the same rank for all $(x, y) \in U_\epsilon(x^0, y^0)$.

Let $L(x, y, \lambda, \mu) = f(x, y) + \lambda^T g(x, y) + \mu^T h(x, y)$ denote the Lagrangian function of problem (13) and let

$$\Lambda(x, y) = \{(\lambda, \mu) \in \mathbb{R}^p \times \mathbb{R}^q : \lambda \geq 0, \lambda^T g(x, y) = 0, \nabla_y L(x, y, \lambda, \mu) = 0\}$$

be the set of Lagrange multipliers at (x, y) .

Theorem 3 (Dempe, 2002) Consider the problem (13) at $x = x^0 \in \mathbb{R}^n$ and let the assumptions (MFCQ), (SSOC), and (CRCQ) be satisfied at (x^0, y^0) with y^0 being a unique local optimal solution. Then, there exists a unique local optimal solution function $y(\cdot)$ that is locally Lipschitz continuous and directionally differentiable at $x = x^0$. The directional derivative in direction r coincides with the unique optimal solution of the following quadratic programming problem

$$\begin{aligned} & \min_d 0.5 d^T \nabla_{yy}^2 L(x^0, y^0, \lambda^0, \mu^0) d + d^T \nabla_{xy}^2 L(x^0, y^0, \lambda^0, \mu^0) r, \\ & \text{s.t. } \nabla_y g_i(x^0, y^0) d + \nabla_x g_i(x^0, y^0) r \begin{cases} = 0, & \text{if } i \in J(\lambda^0), \\ \leq 0, & \text{if } i \in I(x^0, y^0) \setminus J(\lambda^0), \end{cases} \\ & \nabla_y h_j(x^0, y^0) d + \nabla_x h_j(x^0, y^0) r = 0 \text{ for all } j = 1, \dots, q, \end{aligned}$$

for any $(\lambda^0, \mu^0) \in \Lambda(x^0, y^0)$ that solve

$$\max_{(\lambda, \mu) \in \Lambda(x^0, y^0)} \nabla_x L(x^0, y^0, \lambda, \mu).$$

A.2 Proof of Proposition 1

To show the existence of a global optimal solution, it is enough to verify that assumptions (C) and (MFCQ) are satisfied.

Let $g(V, W) = -W$ and $h(V, W) = \sum_{j=1}^J W_j - 1$ denote the constraints in the lower-level problem. Clearly, the set $\{(V, W) \in \mathbb{R}^{K \times K} \times \mathbb{R}^J : g(V, W) \leq 0, h(V, W) = 0\}$ is non-empty and compact. Therefore, condition (C) holds.

To check (MFCQ), let $(V_0, W_0) \in \mathcal{V} \times \mathcal{W}$ and define

$$I_0 = \{j : g_j(V_0, W_0) = -W_j = 0\}.$$

If $W_0 > 0$, we have $I_0 = \emptyset$ and (MFCQ) holds trivially. If there exists at least some index j such that $W_{0j} = 0$, we need to check the gradient conditions. Let $d \in \mathbb{R}^J$ be a candidate direction. From the inequality constraints we have that $\nabla_w g(V, W)d = -d < 0$, which means that for every $j \in I_0$, we require $d_j > 0$. When combined with the equality constraint we have that

$$\nabla_w h(V_0, W_0)d = \sum_{j \in I_0} d_j + \sum_{j \in I_0^c} d_j = 0,$$

where $I_0^c = \{j : g_j(V_0, W_0) \neq 0\}$. Since $h(V_0, W_0) = 0$ and all coefficients cannot be zero, the set I_0^c is non-empty. Therefore, we can find d such that (MFCQ) holds. Now the existence of the optimal solution follows from Theorem 1, which concludes the proof.

A.3 Proof of Proposition 3

Note that the convex combination $X_0 \tilde{W}$ is a K -dimensional vector, where each scalar element $X_{0k} W^*$ is a convex combination of predictor $k = 1, \dots, K$. Suppose $X_{0k} \tilde{W} = X_{1k}$ for some arbitrary k , but not necessarily for other predictors. In this case, it is easy to verify that \tilde{W} remains an optimal solution to the reduced single-dimensional problem using V_k such that the loss function of the lower-level problem goes to zero. Since the lower-level loss function cannot be improved, we have $\tilde{W} \in \Psi(V_k)$ and the solution is considered feasible for the bilevel problem (6)–(7). Furthermore, if the original solution was bilevel optimal, then also the other solution (V_k, \tilde{W}) remains optimal, since the upper-level objective value depends only on \tilde{W} . This concludes the proof.

A.4 Proof of Proposition 4

Given that the assumptions of Theorem 2 are satisfied, the solution set mapping Ψ_{ε_k} of the regularized lower-level problem (8) is upper semi-continuous. That is, for each sequence $\{(V^k, W^k, \varepsilon_k)\}_{k=1}^\infty$ with $\lim_{k \rightarrow \infty} V^k = \bar{V}$, $\lim_{k \rightarrow \infty} \varepsilon_k = 0+$ and $W^k \in \Psi_{\varepsilon_k}(V^k)$ for all k , each accumulation point of the sequence $\{W^k\}_{k=1}^\infty$ is an optimal solution to the lower-level problem, i.e. the accumulation points belong to $\Psi_0(\bar{V}) = \Psi(\bar{V})$. Then, by continuity of L_V the first assertion follows.

To show the second assertion it is enough to verify that the regularized lower-level problem meets the assumptions of Theorem 3. This is easy to check because the requirement that $Y_0^T Y_0$ is positive definite means that $\nabla_{ww} L_V(V, W)$ is positive definite at each $(V, W) \in \mathcal{V} \times \mathcal{W}$, which means that (SSOC) is satisfied at all feasible points. As a result, Theorem 3 implies that the set $\Psi_{\varepsilon_k}(V^k) = \{W^k(V^k)\}$ is a singleton and the optimal solution function $W^k(V^k)$ is uniquely defined and directionally differentiable at each $\varepsilon_k > 0$. The remaining part of the claim follows from the inequality

$$L_W(V^k, W^k(V^k)) \geq \min_{W \in \mathcal{W}} L_W(V^k, W)$$

that holds due to feasibility. As a result, we have that

$$L_V(V^k, W^k(V^k)) \leq \min_W \{L_V(V^k, W) : W \in \Psi(V^k)\},$$

which then implies the last assertion for every fixed $V^k \in \mathcal{V}$. This concludes the proof.

Appendix B: Implementation of SCM Algorithm

B.1 Descent Algorithm Based on Tykhonov Regularization

Based on Proposition 4, the original synthetic control problem can be solved by considering a sequence of single-level problems

$$\min_V \{L_{\varepsilon_k}(V) : V \in \mathcal{V}\} \text{ for } \varepsilon_k \rightarrow 0+, \tag{15}$$

where the implicitly defined objective function $L_{\varepsilon_k}(V) = L_V(V, W_{\varepsilon_k}^*(V))$ is directionally differentiable with respect to V . In the literature on bilevel programming, such an approach is commonly referred to as Tykhonov regularization (Dempe, 2002). This approach is not often available because of the strictness of (SSOC) and (CRCQ) conditions. However, when these criteria are satisfied, they enable the use of algorithms that are essentially similar to gradient descent.

Let $E\Lambda(V, W)$ be the vertex set of lower-level Lagrange multipliers corresponding to point (V, W) ,

$$\Lambda(V, W) = \{(\lambda, \mu) : \lambda \geq 0, \lambda^\top g(V, W) = 0, \nabla_w \mathcal{L}(V, W, \lambda, \mu) = 0\},$$

where $\mathcal{L}(V, W, \lambda, \mu) = L_w^\varepsilon(V, W) + \lambda^\top g(V, W) + \mu^\top h(V, W)$ denotes the Lagrangian function for the regularized lower-level problem. Under (MFCQ) condition, the set $\Lambda(V, W)$ is known to be a non-empty, convex and compact polyhedron. Here functions $g(V, W)$ and $h(V, W)$ denote the vector of lower-level inequality constraints and the equality constraint, respectively.

For a fixed vertex $(\lambda^0, \mu^0) \in \Lambda(V^0, W^0)$ at a point $(V, W) = (V^0, W^0)$, we write $\mathcal{I}(\lambda^0)$ to denote the family of all index sets

$$I \subset I(V^0, W^0) := \{i : g_i(V^0, W^0) = 0\}$$

that satisfy the following two conditions:

- (C1) There is $(\lambda, \mu) \in E\Lambda(V^0, W^0)$ such that $J(\lambda) := \{i : \lambda_i > 0\} \subset I \subset I(V^0, W^0)$.
- (C2) The gradients $\{\nabla_w g_i(V^0, W^0) : i \in I\} \cup \{\nabla_w h(V^0, W^0)\}$ are linearly independent.

Following Dempe (2002), the solution algorithm, which is essentially an adaptation of gradient descent, can be outlined as follows:

procedure TYKHONOV-DESCENT:

Input: Synthetic control problem (6)–(7).

Output: A Bouligand stationary solution.

Step 1: Select $V^0 \in \mathcal{V}$, set $k = 0$, choose $\epsilon, \delta \in (0, 1)$, a small $\epsilon' > 0$, a sufficiently small $\kappa > 0$, and a $w < 0$.

Step 2a: Choose (K^k, λ^k, μ^k) with

$$(\lambda^k, \mu^k) \in E\Lambda(W_{\epsilon_k}^*(V^k), V^k) \text{ and } K^k \in \mathcal{I}(\lambda^k)$$

Compute an optimal solution $(d^k, r^k, \gamma^k, \eta^k, s^k)$ for problem (16).

If $s^k < w$ then go to Step 3.

If $s^k \geq w$ and not all possible samples (λ^k, μ^k, K^k) are tried, then continue with Step 2a.

If all (λ^k, μ^k, K^k) have been tried, set $w = w/2$.

If $|w| < \epsilon'$, go to Step 2b, otherwise continue with Step 2a.

Step 2b: Choose (K^k, λ^k, μ^k) satisfying

$$K^k \subset I_\kappa(W_{\epsilon_k}^*(V^k), V^k) \text{ and (C2) as well as}$$

$$(\lambda^k, \mu^k) \in \underset{(\lambda, \mu)}{\operatorname{argmin}} \{ \|\nabla_w \mathcal{L}(W_{\epsilon_k}^*(V^k), V^k, \lambda, \mu)\|^2 : \lambda_j = 0, j \notin K^k \}.$$

Here $I_\kappa = \{j : -\kappa \leq g_j(V, W) \leq 0\}$ denotes the set of κ -active lower-level inequalities.

Compute an optimal solution $(d^k, r^k, \gamma^k, \eta^k, s^k)$ for problem (16).

If $s^k < w$, go to Step 3.

If $s^k \geq w$ and not all (λ^k, μ^k, K^k) have been tried, continue with Step 2b.

If all K^k have been tried, then set $w = w/2$. If $|w| < \epsilon'$, then stop.

Step 3: Choose a largest step-size $t^k \in \{\delta, \delta^2, \delta^3, \delta^4, \dots\}$ such that

$$L_{\epsilon_k}(V^k + t^k r^k) \leq L_{\epsilon_k}(V^k) + \epsilon t^k s^k, \quad G(V^k + t^k r^k) \leq 0.$$

If $t^k < \epsilon'$, then drop the actual set K^k and continue searching for a new set K^k in Step 2a or 2b.

Step 4: Set $V^{k+1} = V^k + t^k r^k$, $k = k + 1$.

Step 5: If a stopping criterion is satisfied, i.e. ϵ_k is sufficiently small, then stop. Otherwise, set $\epsilon_{k+1} = \delta \epsilon_k$ and compute $W_{\epsilon_{k+1}}^*(V^{k+1})$ and go to step 2.

end procedure

The directional derivative in Step 2 can be computed using quadratic programming based on Theorem 3 by Dempe (2002). Let $K^k \in \mathcal{I}(\lambda^k)$ be some index set and $v^k = (\lambda^k, \mu^k) \in E\Lambda(z^k)$ be a vertex, where $z^k = (V^k, W^k)$. Then the descent direction r^k is obtained as part of a solution to the following problem:

$$\begin{aligned}
 & \min_{d,r,\gamma,\eta,s} s & (16) \\
 \text{s.t. } & L'_{\epsilon_k}(V^k; r^k) := \nabla_w L_V(z^k)d + \nabla_v F(z^k)r \leq s \\
 & \nabla_v G_i(V^k)r \leq -G_i(V^k) + s, \quad i = 1, \dots, K + 2 \\
 & \nabla_{ww}^2 \mathcal{L}(z^k, v^k)r + \nabla_w^T g(z^k)\gamma + \nabla_w^T h(z^k)\eta = 0 \\
 & \nabla_w g_i(z^k)d + \nabla_v g_i(z^k)r \begin{cases} = 0, & i \in K^k \\ \leq -g_i(z^k) + s, & i \notin K^k \end{cases} \\
 & \nabla_w h(z^k)d + \nabla_v h(z^k)r = 0 \\
 & \lambda_i + \gamma_i + s \geq 0, \quad i \in K^k, \quad \gamma_i = 0, \quad i \notin K^k, \quad \|r\| \leq 1.
 \end{aligned}$$

When the problem has a feasible solution $(d^k, r^k, \gamma^k, \eta^k, s^k)$ such that the objective value is negative, $s^k < 0$, for some index set K^k and vertex v^k , then the point (V^k, W^k) is not locally optimal. This means that there exists a direction r^k for which the directional derivative of L_{ϵ_k} is negative at V^k .

When parametrizing the algorithm, it is useful to choose the value for ϵ' to be small enough to ensure that Step 3 terminates only if a set K^k is selected in Step 2b such that the problem (16) has a negative optimal value. It is also noteworthy that Step 2b should be considered only when the value of $L_{\epsilon_k}(V^k; r^k)$ is sufficiently small and even then only for small κ . Otherwise, there is a risk of increasing numerical effort substantially. For discussion on the convergence of this kind of algorithm to a Bouligand stationary point, we refer to Dempe and Schmidt (1996).

B.2 Algorithm Based on KKT Approximations

The use of KKT reformulations has been a common practice when solving bilevel problems. Unfortunately, this has turned out to be far more difficult than anticipated. Quite commonly, the local optimal solutions obtained by solving KKT reformulated problems do not correspond to the local optimal solutions of the original bilevel problem. While the KKT reformulations are equivalent to the original problem in terms of global optimal solutions, the equivalence is lost when numerical algorithms need to be used. Since KKT reformulations typically lead to a nonconvex optimization problem, the solution algorithms tend to find only stationary or local optimal solutions, which may not correspond to the solutions of the original problem.

Fortunately, there is still some good news left when it comes to the use of KKT conditions in practice. In their recent paper, Dempe and Franke (2019) suggest a numerically stable approach for handling optimistic bilevel problems with convex lower-level problems. The idea is based on a clever approximation of the KKT transformation which enables us to use general solution algorithms for

non-convex optimization problems to approximate the local optimal solution of the original bilevel optimization problem.

Now instead of considering the classical KKT reformulation of the problem, the idea developed in the paper by Dempe and Franke (2019) is to construct perturbed problems that approximate the original formulation. Let \mathcal{L} denote the Lagrangian corresponding to the lower-level problem,

$$\mathcal{L}_\varepsilon(V, W, \lambda) = L_W^\varepsilon(V, W) + \lambda^\top g(V, W).$$

We then solve a sequence of perturbed problems

$$\begin{aligned} \min_{V, W, \lambda} \quad & L_V(V, W) \\ & G(V) \leq 0 \\ & \|\nabla_w \mathcal{L}_\varepsilon(V, W, \lambda)\| \leq e_1 \\ & g(V, W) \leq 0 \\ & \lambda \geq 0, \\ & -\lambda_i g_i(V, W) \leq e_2, \quad i = 1, \dots, J + 2, \end{aligned} \tag{17}$$

for $(e_1, e_2) \rightarrow 0+$ and $\varepsilon \rightarrow 0+$. Here, the norm $\|\cdot\|$ can be chosen to be for instance the Chebyshev norm $\|a\|_\infty = \max_{i=1, \dots, n} |a_i|$ or the usual Euclidean norm $\|a\|_2 = \sqrt{\sum_{i=1}^n a_i^2}$. The function G is defined such that it matches the definition of set $\mathcal{V} = \{V : G(V) \leq 0\}$ in (5). Similarly, g represents the lower-level constraints such that $\mathcal{W} = \{W : g(V, W) \leq 0\}$ corresponds to (4).

Earlier, a similar approach of using a sequence of perturbed problems to solve bilevel problems has also been considered by Mersha and Dempe (2011), who suggested a specifically tailored algorithm to solve the problem. Later, however, Dempe and Franke (2019) have shown that the assumptions made earlier have been too restrictive and the sequence of perturbed problems can actually be solved by an arbitrary algorithm.

Author Contributions All authors contributed to the study. PM Writing—original draft, Conceptualization, Methodology. JE Conceptualization, Writing—review & editing. XZ Writing—original draft, Writing—review & editing, Conceptualization, Visualization, Software, Formal analysis. TK Writing—original draft, Writing—review & editing, Conceptualization, Methodology, Formal analysis. All authors read and approved the final manuscript.

Funding The authors declare that no funds, grants, or other support were received during the preparation of this manuscript.

Data Availability The data that support the findings of this study are openly available in the Matlab implementation of the *Synth* package: <https://web.stanford.edu/~jhain/synthpage.html> (accessed 24 April 2023).

Declarations

Competing interests The authors have no relevant financial or non-financial interests to disclose.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Abadie, A. (2020). Statistical nonsignificance in empirical economics. *American Economic Review: Insights*, 2(2), 193–208.
- Abadie, A. (2021). Using synthetic controls: Feasibility, data requirements, and methodological aspects. *Journal of Economic Literature*, 59(2), 391–425.
- Abadie, A., Diamond, A., & Hainmueller, J. (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., Diamond, A., & Hainmueller, J. (2011). Synth: An R package for synthetic control methods in comparative case studies. *Journal of Statistical Software*, 42(13), 1–17.
- Abadie, A., Diamond, A., & Hainmueller, J. (2015). Comparative politics and the synthetic control method. *American Journal of Political Science*, 59(2), 495–510.
- Abadie, A., & Gardeazabal, J. (2003). The economic costs of conflict: A case study of the Basque country. *American Economic Review*, 93(1), 113–132.
- Acemoglu, D., Johnson, S., Kermani, A., Kwak, J., & Mitton, T. (2016). The value of connections in turbulent times: Evidence from the United States. *Journal of Financial Economics*, 121(2), 368–391.
- Albalade, D., Bel, G., & Mazaira-Font, F. A. (2021). Decoupling synthetic control methods to ensure stability, accuracy and meaningfulness. *SERIEs*.
- Alfano, V., Ercolano, S., & Cicatiello, L. (2021). School openings and the COVID-19 outbreak in Italy. A provincial-level analysis using the synthetic control method. *Health Policy*, 125(9), 1200–1207.
- Amjad, M., Shah, D., & Shen, D. (2018). Robust synthetic control. *Journal of Machine Learning Research*, 19(1), 802–852.
- Arkhangelsky, B. D., Athey, S., Hirshberg, D. A., Imbens, G. W., & Wager, S. (2021). Synthetic difference-in-differences. *American Economic Review*, 111(12), 4088–4118.
- Athey, S., & Imbens, G. W. (2017). The state of applied econometrics: Causality and policy evaluation. *Journal of Economic Perspectives*, 31(2), 3–32.
- Bank, B., Guddat, J., Klatte, D., Kummer, B., & Tammer, K. (1982). *Non-linear parametric optimization*. Berlin: Akademie-Verlag.
- Bayer, P., & Aklin, M. (2020). The European Union Emissions Trading System reduced CO2 emissions despite low prices. *Proceedings of the National Academy of Sciences of the United States of America*, 117(16), 8804–8812.
- Becker, M., & Klößner, S. (2017). Estimating the economic costs of organized crime by synthetic control methods. *Journal of Applied Econometrics*, 32(7), 1367–1369.
- Becker, M., & Klößner, S. (2018). Fast and reliable computation of generalized synthetic controls. *Econometrics and Statistics*, 5, 1–19.
- Becker, M., Klößner, S., & Pfeifer, G. (2018). Cross-validating synthetic controls. *Economics Bulletin*, 38(1), 603–609.
- Ben-Michael, E., Feller, A., & Rothstein, J. (2021). The augmented synthetic control method. *Journal of the American Statistical Association*, 116(536), 1789–1803.
- Bonander, C., Humphreys, D., & Esposti, M. D. (2021). Synthetic control methods for the evaluation of single-unit interventions in epidemiology: A tutorial. *American Journal of Epidemiology*, 190(12), 2700–2711.

- Cavallo, E., Galiani, S., Noy, I., & Pantano, J. (2013). Catastrophic natural disasters and economic growth. *Review of Economics and Statistics*, 95(5), 1549–1561.
- Charnes, A. (1952). Optimality and degeneracy in linear programming. *Econometrica*, 20(2), 160–170.
- Cole, M. A., Elliott, R. J., & Liu, B. (2020). The impact of the Wuhan Covid-19 lockdown on air pollution and health: A machine learning and augmented synthetic control approach. *Environmental and Resource Economics*, 76(4), 553–580.
- Dempe, S. (2002). *Foundations of bilevel programming*. Dordrecht: Kluwer.
- Dempe, S., & Franke, S. (2019). Solution of bilevel optimization problems using the KKT approach. *Optimization*, 68(8), 1471–1489.
- Dempe, S., & Schmidt, H. (1996). On an algorithm solving two-level programming problems with nonunique lower level solutions. *Computational Optimization and Applications*, 6(3), 227–249.
- Douchchenko, N., & Imbens, G. W. (2017). Balancing, regression, difference-in-differences and synthetic control methods: A synthesis. arXiv preprint [arXiv:1610.07748](https://arxiv.org/abs/1610.07748).
- Gobillon, L., & Magnac, T. (2016). Regional policy evaluation: Interactive fixed effects and synthetic controls. *Review of Economics and Statistics*, 98(3), 535–551.
- Hansen, P., Jaumard, B., & Savard, G. (1992). New branch-and-bound rules for linear bilevel programming. *SIAM Journal on Scientific and Statistical Computing*, 13(5), 1194–1217.
- Kaul, A., Klößner, S., Pfeifer, G., & Schieler, M. (2022). Standard synthetic control methods: The case of using all preintervention outcomes together with covariates. *Journal of Business and Economic Statistics*, 40(3), 1362–1376.
- Kleven, H. J., Landais, C., & Saez, E. (2013). Taxation and international migration of superstars: Evidence from the European football market. *American Economic Review*, 103(5), 1892–1924.
- Klößner, S., & Pfeifer, G. (2015). Synthesizing Cash for Clunkers: Stabilizing the Car Market, Hurting the Environment. Beiträge zur Jahrestagung des Vereins für Socialpolitik 2015: Ökonomische Entwicklung - Theorie und Politik - Session: Automobiles and the Environment, F13-V1.
- Klößner, S., Kaul, A., Pfeifer, G., & Schieler, M. (2018). Comparative politics and the synthetic control method revisited: A note on Abadie et al. (2015). *Swiss Journal of Economics and Statistics*, 154(1), 11.
- Kuosmanen, T., Zhou, X., Eskelinen, J., & Malo, P. (2021). Design flaw of the synthetic control method. MPRA Paper, 106390.
- Lang, D., Esbenshade, L., & Willer, R. (2022). Did Ohio's vaccine lottery increase vaccination rates? A pre-registered, synthetic control study. *Journal of Experimental Political Science*, 2022, 1–19.
- Mersha, A. G., & Dempe, S. (2011). Direct search algorithm for bilevel programming problems. *Computational Optimization and Applications*, 49(1), 1–15.
- Mills, M. C., & Rüttenauer, T. (2022). The effect of mandatory COVID-19 certificates on vaccine uptake: Synthetic-control modelling of six countries. *The Lancet Public Health*, 7(1), e15–e22.
- Mitze, T., Kosfeld, R., Rode, J., & Walde, K. (2020). Face masks considerably reduce COVID-19 cases in Germany. *Proceedings of the National Academy of Sciences of the United States of America*, 117(51), 32293–32301.
- Sehgal, N. K. (2021). Impact of vax-a-million lottery on COVID-19 vaccination rates in Ohio. *American Journal of Medicine*, 134(11), 1424–1426.
- Sinha, A., Malo, P., & Deb, K. (2013). Efficient evolutionary algorithm for single-objective bilevel optimization. arXiv preprint [arXiv:1303.3901](https://arxiv.org/abs/1303.3901).
- Vicente, L., Savard, G., & Júdice, J. (1994). Descent approaches for quadratic bilevel programming. *Journal of Optimization Theory and Applications*, 81(2), 379–399.
- Xin, M., Shalaby, A., Feng, S., & Zhao, H. (2021). Impacts of COVID-19 on urban rail transit ridership using the Synthetic Control Method. *Transport Policy*, 111(June), 1–16.
- Xu, Y. (2017). Generalized synthetic control method: Causal inference with interactive fixed effects models. *Political Analysis*, 25(1), 57–76.