

State-Building through Public Land Disposal? An Application of Matrix Completion for Counterfactual Prediction

Jason Poulos[†]

Dept. of Health Care Policy
Harvard Medical School

December 29, 2023

Abstract

This paper examines how homestead policies, which opened vast frontier lands for settlement, influenced the development of American frontier states. It uses a treatment propensity-weighted matrix completion model to estimate the counterfactual size of these states without homesteading. In simulation studies, the method shows lower bias and variance than other estimators, particularly in higher complexity scenarios. The empirical analysis reveals that homestead policies significantly and persistently reduced state government expenditure and revenue. These findings align with continuous difference-in-differences estimates using 1.46 million land patent records. This study's extension of the matrix completion method to include propensity score weighting for causal effect estimation in panel data, especially in staggered treatment contexts, enhances policy evaluation by improving the precision of long-term policy impact assessments.

Keywords: Causal inference; Difference-in-differences; Matrix completion; State size; Synthetic controls; Panel data.

[†]*Corresponding Author.* jvpoulos@bwh.harvard.edu. Division of Endocrinology, Brigham and Women's Hospital, 221 Longwood Avenue, Boston, MA 02115.

1. Introduction

The exploration of state development patterns over time and across regions is a growing area of interest for social scientists. A key contribution in this field comes from Bense (1990, p. 164), who emphasizes the significant role of mid-nineteenth century homestead policies — federal laws aimed at transferring public land to private individuals — in shaping the developmental trajectory of the United States. Additionally, Murtazashvili (2013, p. 250) and Frymer (2017, p. 12) suggest that these policies not only facilitated land distribution but also enhanced the federal government’s bureaucratic capacity to manage public lands and secure future revenue streams. This paper examines how homestead policies impacted the size of state governments, which is closely related to state capacity, or the ability of governments to finance and implement policies (Besley and Persson, 2010).

Homesteading is expected to expand the size of state governments by increasing the land values and tax revenue of sparsely populated frontier states. The expansion in state size was historically evident in the adoption of compulsory primary education laws and public education investments by frontier state governments, as a strategy to attract homesteaders (Engerman and Sokoloff, 2005). Contrary to this expectation, I provide evidence that homesteads authorized under the Homestead Act (HSA) of 1862 and the Southern Homestead Act (SHA) of 1866, which opened for settlement hundreds of millions of acres of land for homesteading, significantly reduced the size of frontier state governments over the long-run. The finding that the homestead acts limited the size of frontier governments aligns with Mattheis and Raz’s (2021) findings that regions impacted by the HSA experienced a slower transition from agriculture to other economic sectors and lower housing values. The paper further investigates land inequality as a possible causal mechanism underlying the relationship between homestead policies and state capacity, considering that median voter-based theories of inequality and redistribution predict inequality increases the size of governments, and show that exposure to homesteads decreased land inequality

over time.

The paper makes a methodological contribution in extending the matrix completion method (Athey et al., 2021) for estimating the causal effects of policy interventions in panel data, by weighting the loss function with estimated unit- and time-varying treatment probabilities (i.e., the propensity score) to correct for imbalances in the covariate distributions between the factual and counterfactual values. This extension, which was proposed by Athey et al. (2021) but has not been implemented, places more emphasis on the loss for factual unit-time values that are most similar to the counterfactual values in terms of pre-treatment covariates. The covariates used in the application control for selective migration to more agriculturally productive land, and for selection bias arising from differences in access to frontier lands.

A standard method for causal inference with panel data is difference-in-differences (DID), which relies on the parallel trends assumption: in the absence of treatment, the average outcomes of treated and control units would have followed parallel paths. Under parallel trends, DID identifies causal effects by contrasting the change in outcomes pre- and post-treatment, between the treated and control groups. However, the parallel trends assumption is generally invalid in the presence of unobserved time-varying confounders. DID has been extended to staggered treatment implementation settings, where the time of initial treatment varies among multiple treated units (Callaway and Sant’Anna, 2020; Goodman-Bacon, 2021; Athey and Imbens, 2021).

Another popular method of handling unobserved time-varying confounders in panel data is the synthetic control method (SCM; Abadie et al., 2010). The method constructs a convex combination of control units that are similar to a single treated unit in terms of pre-treatment outcomes or covariates, to help balance unobserved time-varying confounding between treatment and control groups. The SCM estimator assumes there is a stable convex combination of the control units that absorbs all time-varying unobserved confounding. The convexity restriction is equivalent to imposing a restriction of linear dependence between

factor loadings in the context of matrix completion or latent factor models (Gobillon and Magnac, 2016; Xu, 2017; Xiong and Pelger, 2020; Bai and Ng, 2021). The SCM can be generalized to settings with multiple treated units (Doudchenko and Imbens, 2016) and staggered treatments (Ben-Michael et al., 2019), and to include features of DID estimation (Ben-Michael et al., 2018; Arkhangelsky et al., 2021) or matrix factorization (Amjad et al., 2018; Agarwal et al., 2021; Fan et al., 2021).

Similar to latent factor models, matrix completion attempts to model unobserved time-varying confounders by decomposing the factual outcomes into matrices of latent factors (i.e., time-varying coefficients) and factor loadings (i.e., unit-specific intercepts). The counterfactual values are then imputed using the estimated factors and loadings. Matrix completion and latent factor models avoid imposing convexity constraints on the factor loadings like the SCM, and typically use matrix norm regularization or factorization to produce a low-dimensional representation of the factual outcomes, which improves generalizability. Matrix completion offers distinct advantages over latent factor models: first, it does not require fixing the rank (i.e., number of unobserved factors) of the underlying data; second, it is suitable in staggered treatment implementation settings, even when few control units are available; third, it uses all factual data to estimate unobserved factors, while latent factor models use only the pre-treatment data.

The structure of this paper is organized as follows: Section 2 provides an overview the historical context of homestead policies and their relationship to state size and land inequality. Section 3 details the matrix completion method for obtaining the causal estimands of interest and reports the results of simulation studies to evaluate the proposed method. Section 4 describes the data sources used for the application and potential sources of bias in the analysis. Section 5 presents estimates of the long-run impacts of homestead policies on state size using the matrix completion estimator and reports the results of a “no-treatment evaluation” to verify the consistency of the estimator. Section 6 reports DID estimates of the effect of homesteads on state size and land inequality. The final section discusses the

study’s findings, focusing on the long-term negative effects of homestead policies on state finances and the role of land inequality in state capacity. It also emphasizes the study’s methodological advancement in policy evaluation through the matrix completion method with propensity score weighting.

2. Historical background

The view that the western frontier had long-lasting impacts on the evolution of democratic institutions can be traced to Turner (1956). Turner’s frontier thesis posits that homestead policies acted as a “safety valve” for relieving pressure from congested urban labor markets in eastern states. The view of the frontier as a safety valve has been further explored by Ferrie (1997), who finds evidence in a linked census sample of substantial migration to the frontier by unskilled workers and considerable gains in wealth for these migrant workers. Bazzi et al. (2020) expand on this demographic profile, showing that frontier settlers, often illiterate and foreign-born, possessed a distinct individualism suited for the challenging frontier life. The historical experience of the frontier is reflected in modern times through lower property tax rates in counties with a longer frontier history and a prevailing sentiment among residents against taxation and redistribution.

Homestead policies not only offered greater economic opportunities to eastern migrants, but also the sparse population on the western frontier meant that state and local governments competed with each other to attract migrants in order to lower local labor costs and to increase land values and tax revenue. Frontier governments offered migrants broad access to cheap land and property rights, unrestricted voting rights, and more generous provision of schooling and other public goods (Engerman and Sokoloff, 2005). Consistent with this view, Poulos and Zeng (2021) estimate that the long-run impact of homestead policies on public school spending is equivalent to 2.5% of the total per-capita public school expenditures in 1929.

García-Jimeno and Robinson (2008) test the frontier thesis in a global context and conclude that the economic effect of the frontier depends on the quality of political institutions at the time of frontier expansion. Frontier expansion promotes equitable outcomes only when societies are initially democratic. When institutional quality is weak, the existence of frontier land can yield worse developmental outcomes because non-democratic political elites can monopolize frontier lands.

2.1. Homestead policies

The 1862 HSA opened up hundreds of millions of acres of western public land for settlement. Any adult household head — including women, immigrants who had applied for citizenship, and freed slaves following the passage of the Fourteenth Amendment— could apply for a homestead grant of 160 acres of land, provided that they live and make improvements on the land for five years and pay a \$10 filing fee. Under the HSA, the bulk of newly surveyed land on the western frontier was reserved for homesteads, although the law did not end sales of public land. The explicit goal of the HSA was to liberalize the homesteading requirements set by the Preemption Act of 1841, which permitted individuals already inhabiting public land to purchase up to 160 acres at \$1.25 per acre before the land was put up for sale. The implicit goal was to promote rapid settlement of the western frontier (Allen, 1991), and to reduce the federal government’s costs of defending contested frontiers (Frymer, 2014).

In the pre-Reconstruction South, public land was not open to homestead but rather unrestricted cash entry, which permitted the direct sale of public land to private individuals of 80 acres or more for at least \$1.25 an acre. The 1866 SHA restricted cash entry and reserved for homesteading over 46 million acres of public land, or about one-third of the total land area in the five southern public land states (Lanza, 1999, p. 13). The Bureau of Land Management (BLM) classified land disposed under the SHA under the same authority as land disposed by the HSA, since the SHA amended the HSA and dictated that public lands be disposed under the stipulations of the HSA (Hoffnagle, 1970).

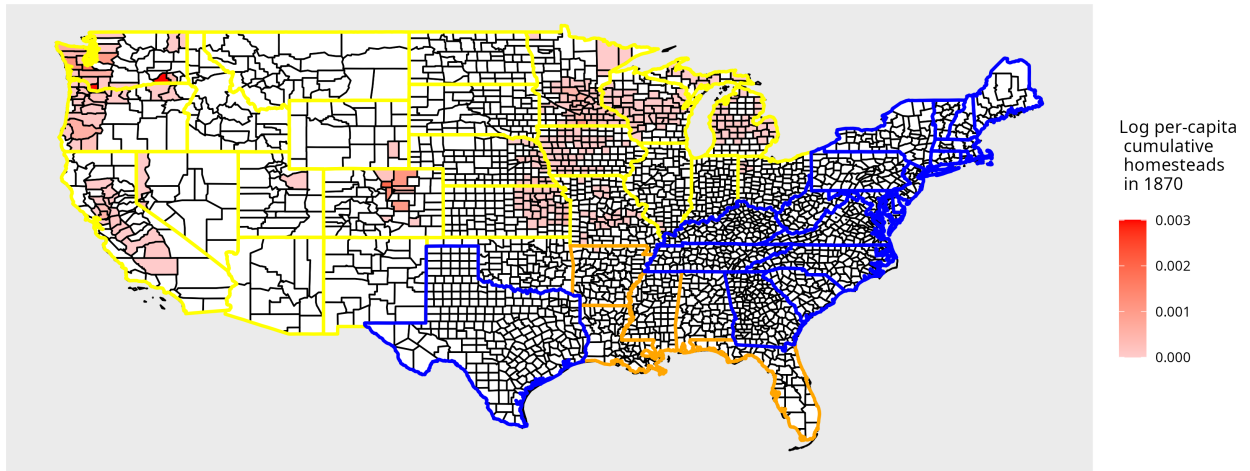
2.1.1. Public and state-land states

Public land states (PLS) are states that were crafted from the public domain, and where the federal government has the primary authority to distribute public land (Murtazashvili, 2013, p. 4). In the South, these states include Alabama, Arkansas, Florida, Louisiana, and Mississippi. Western PLS include the 25 states that comprise the Midwestern, Southwestern, and Western U.S. (except Hawaii). State-land states, which include the original 13 states, Kentucky, Maine, Tennessee, Texas, Vermont, and West Virginia, were not open to homesteading because the state governments had primary authority to distribute public land. The maps in Figure 1 reflect the impact of these policies, indicating a significant increase in log per-capita cumulative homesteads in Southern and Western PLS by 1900.

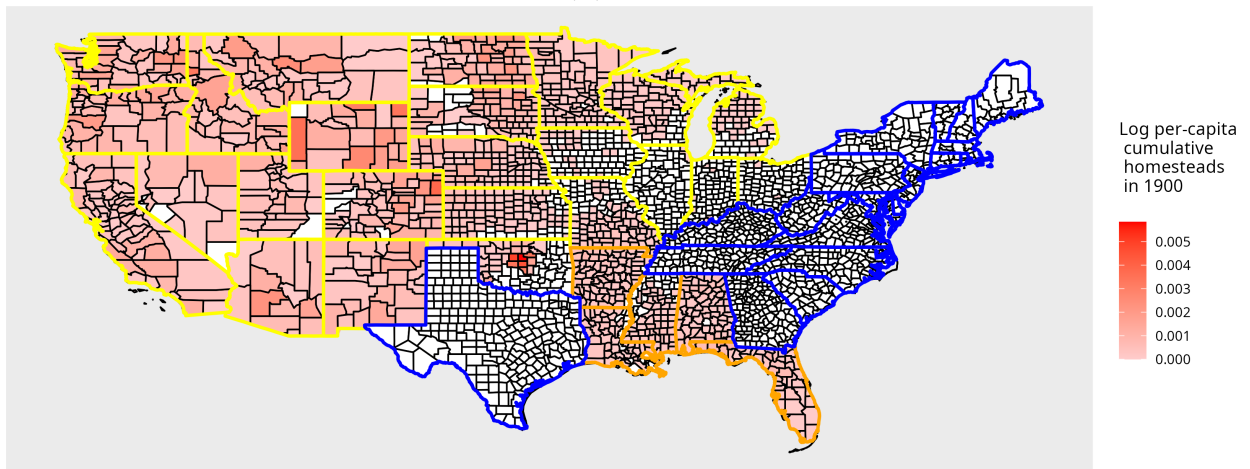
2.1.2. Challenges and speculation in homesteading

There were substantial barriers to entry to homesteading, and homesteaders took on enormous risk in the five years required to file a homestead patent. One of the most significant obstacles to entry was the need for capital to build a successful farming operation: contemporary writers estimated that potential homesteaders required \$600 to \$1000 to start a farm (Deverell, 1988). The high cost associated with starting and maintaining a farm casts doubt on the safety valve hypothesis (Danhof, 1941), and the effectiveness of the land policies such as the HSA as a wealth-building tool was limited by the binding capital constraints faced by small farmers (Gates and Bogue, 1996, p. 35). Poulos et al. (2023) further contributes to this discussion by examining the 1901 Oklahoma Land Lottery, demonstrating gender differences in leveraging lottery wealth for land purchases and homestead patents, with female winners more effectively using lottery wealth to overcome liquidity constraints in entrepreneurial activities.

Homesteading was a risky venture: over the period of 1910 to 1919, out of 604,092 homestead entries in the U.S., totaling over 128 million acres, only 384,954 (63.7%) resulted in



(A) 1870



(B) 1900

Figure 1: Log per-capita cumulative homesteads in 1870 and 1900, overlaid on 1911 county borders (Long, 1995). White-colored counties have no homestead entries. States bordered in blue (□) are state land states; yellow (□) denote western public land states; orange (□) denote southern public land states.

successful patents (Shanks, 2005). At least part of the discrepancy between homestead entries and filings, however, may be explained by fraudulent filings. Speculators and corporations engaged in the practice of paying individual to stake a claim in a homestead, with no intention of completing the patent, in order to extract resources from the land (Gates, 1936). In the South, these “dummy entrymen” were used by timber and mining companies to extract resources while the cash entry restriction of the SHA was in effect. When the restriction was removed, there was no need for fraudulent filings because the larger companies could buy land in unlimited amounts at a nominal price (Gates, 1940, 1979). The same pattern of fraudulent filings existed in the West, where Murtazashvili (2013, p. 216) argues that speculators benefited disproportionately from public land policies because the economic balance of power tilted toward the wealthy. Gates (1942) characterizes western speculators who bought land in bulk prior to the 1889 restriction as being influential in state and local governments, resistant to paying taxes, and opposed to government spending.

2.2. Land inequality as a causal mechanism

Inequality is a potential causal mechanism underlying the relationship between homesteads and state size. Median voter-based theories that assume parity in the political influence of voters predict a positive relationship between inequality and the size of governments (Meltzer and Richard, 1981). In settings with high inequality, the median voter is poorer than the average voter, which in turn increases demand for redistribution in majority-rule elections.

However, models that allow for economic differences in the political influence of voters predict a nonlinear or inverse relationship between inequality and government size. In Benabou’s (2000) model, for instance, the pivotal voter is wealthier than the median and has the power to block redistribution as inequality increases. But when inequality is too high, the poor can impose redistribution on elites through majority voting (Perotti, 1993; Saint-Paul and Verdier, 1993). In Besley and Persson’s (2009) framework, for example, greater

economic power of the ruling class reduces government spending and investments in state capacity. Similarly, Galor et al. (2009) propose a model where wealthy landowners block education reforms because education favors industrial labor productivity and decreases the value in farm rents. Inequality in this context can be thought of as a proxy for the amount of *de facto* political influence elites have to block reforms and limit the size of states (Acemoglu and Robinson, 2008).

To test whether homesteads affected future land inequality in frontier counties, I calculate a commonly-used measure of land inequality based on the Gini coefficient of census farm sizes, adjusted by the ratio of farms to adult males, a measure proposed by Vollrath (2013). Gini-based land inequality measures are commonly used as proxy for the *de facto* bargaining power of landed elites (e.g., Boix, 2003; Ziblatt, 2008; Ansell and Samuels, 2015).

Figure A1 in the online Appendix plots the results from bivariate regression models of land inequality and state government finances during the period of 1860 to 1950, demonstrating a positive relationship among groups of PLS and state-land states. This associational evidence is consistent with the predictions of the Meltzer and Richard model; however, it contrasts with recent empirical studies that establish a negative relationship based on cross-sectional analyses. Ramcharan (2010), for instance, finds an inverse relationship between land inequality and county-level property tax revenue in 1890. The authors find that the negative relationship is especially large in rural counties, where landownership tends to be more concentrated. Vollrath (2013) establish a negative relationship between land inequality and local property tax revenue in northern rural counties in 1890. The present findings, in contrast, are based on state-level expenditure and revenue panel data.

3. Matrix completion for counterfactual prediction

This paper applies the matrix completion method proposed by Athey et al. (2021) to predict counterfactual outcomes, and extends the method by propensity-weighting the loss

function to correct for imbalances in the factual covariate distributions between treatment and control groups.

3.1. Setup and notation

We consider a sample of $i \in \{1, \dots, N\}$ units, each observed in $t \in \{1, \dots, T\}$ time periods. Following the notation of Athey and Imbens (2021), let \mathbf{a} be a length- N vector, where $a_i \in \{1, \dots, T, \infty\}$ indexes the time of initial treatment, and $a_i = \infty$ denotes control units. If a unit enters treatment during the panel ($a_i \neq \infty$), it remains treated for the remainder of the panel. There is a nonzero number of control units, $N_C = N - \sum_i \mathbb{1}_{a_i = \infty}$, with $\mathbb{1}(\cdot)$ denoting the indicator function, and a nonzero number treated units $N_T = N - N_C = \sum_i \mathbb{1}_{a_i \neq \infty}$. Let the values of the treatment indicator $W_{it} \in \{0, 1\}$ be $W_{it} = 0$ for the control units in all time periods and $W_{it} = 1$ for the treated units when $t \geq a_i$. Let \mathcal{O} denote the set of factual outcome values; i.e., the values for which $W_{it} = 0$.

Under the Neyman-Rubin potential outcomes framework (Neyman, 1923; Rubin, 1990), for each unit i and time t , there exist potential outcomes $Y(a_i)_{it}$. The fundamental problem is that we can only observe a single potential outcome for each unit-time observation: $Y(a_i)_{it}$ is observed for treated units when entering treatment, and $Y(\infty)_{it}$ is observed for the control units in all time periods. The potential outcomes framework implicitly assumes treatment is well-defined to ensure that each unit has the same number of potential outcomes. It also requires that the potential outcomes of unit i varies with a_i but not with the other values of \mathbf{a} , which is often referred to as the no interference assumption.

There are two additional assumptions are needed to write potential outcomes as a function of \mathbf{a} , which are both made in Athey and Imbens (2021). First, there are no anticipatory effects; i.e., $Y(a_i)_{it} = Y(\infty)_{it}$ for all $a_i > t$. This assumption, which is often implicitly made in panel data studies, assumes that if a unit has not yet entered treatment, the initial treatment time has no causal effect on potential outcomes in the current period. Second, potential outcomes in period t are invariant to how long unit i has been exposed

to treatment; i.e., $Y(a_i)_{it} = Y(1)_{it}$ for all $a_i \leq t$. This assumption does not rule out causal effects of treatment duration on the outcome, but rather rules out causal effects varying by initial treatment time.

3.2. Causal estimands

The causal estimand of interest is the average treatment effect on the treated units (ATT) of entering treatment in a'_i relative to being control ($a_i = \infty$), on the outcome in period t :

$$\tau_{t,\infty a'_i} = \frac{1}{N_T} \sum_{i=1}^{N_T} Y(a'_i)_{it} - Y(\infty)_{it}, \quad \text{for } W_{it} = 1. \quad (1)$$

In the application, I consider $a'_i = \min_{1 \leq i \leq N_T} a_i$, or the year of the earliest homestead entry among the treated units. The ATT averaged over the counterfactual period, which provides a summary measure of the overall treatment effect, is also of interest:

$$\tau_{\infty a'_i} = \frac{1}{T - a'_i + 1} \sum_{t=a'_i}^T \tau_{t,\infty a'_i}. \quad (2)$$

3.3. Estimation

In the application, the outcome of interest is state size, measured by state government spending and revenue. The goal is to estimate the potential outcomes under control for the treated units; i.e., the counterfactual state size of treated units had they not been exposed to treatment. I model the outcome under control as:

$$Y(\infty)_{it} = L_{it} + \gamma_i + \delta_t + \epsilon_{it}, \quad (3)$$

where L_{it} is a typical element in the unknown matrix, $\mathbf{L} = \mathbf{U}\mathbf{V}^\top$, the product of a matrix of factor loadings, $\mathbf{U}_{N \times R}$, and a matrix of factors, $\mathbf{V}_{T \times R}$. While latent factor models

assume the rank, or number of unobserved factors R , is fixed, matrix completion methods assume that the rank of \mathbf{L} is low relative to N and T . The model includes unit-specific fixed effects, $\{\gamma_i\}_{i=1}^N$, and time-specific fixed effects, $\{\delta_t\}_{t=1}^T$, which are meant to capture unobserved confounders not absorbed by the low-rank matrix. The identifying assumption is that the errors ϵ_{it} are conditionally mean zero and independent of a_i , for all values of i and t :

$$\mathbb{E}(\epsilon_{it}|L_{it}, \gamma_i, \delta_t) = \mathbb{E}(\epsilon_{it}|L_{it}, \gamma_i, \delta_t, a_i) = 0. \quad (4)$$

This assumption rules out correlation between the errors and initial treatment time in any period (Ben-Michael et al., 2019). It is analogous to the strict exogeneity assumption made for the estimation of the ATT using latent factor models (Xu, 2017).

Estimating \mathbf{L} involves minimizing the sum of squared errors via nuclear norm regularized least squares:

$$\arg \min_{\mathbf{L}, \gamma, \delta} \left[\frac{1}{|\mathcal{O}|} \sum_{(i,t) \in \mathcal{O}} \frac{\widehat{w}_{it}}{1 - \widehat{w}_{it}} \left(Y(\infty)_{it} - L_{it} - \gamma_i - \delta_t \right)^2 + \lambda_L \|\mathbf{L}\|_{\star} \right], \quad (5)$$

where the nuclear norm, $\|\cdot\|_{\star} = \sum_i \sigma_i(\cdot)$, or sum of singular values, is used to yield a low-rank solution for \mathbf{L} . The value of the hyperparameter λ_L is chosen among 30 possible values by five-fold cross-validation, where in each fold, 80% of the entries in \mathcal{O} are randomly selected to be used for training, while the remaining 20% of entries are used for model validation. The model with λ_L values that yield the lowest root mean squared error averaged over the validation sets is then fit using all entries in \mathcal{O} .

To quantify the propensity score, w_{it} , I model the probability of treatment as:

$$w_{it} = \Pr \left(W_{it} = 1 | Y_{i,1}, \dots, Y_{i,a'_i-1} + X_{ip} \right), \quad 0 < w_{it} < 1, \quad (6)$$

where X_{ip} is a typical element in a matrix of p covariates measured prior to a'_i . In the application, the covariates include state-level measures of racial composition, prevalence of

pre-emancipation slavery, average farm sizes and values, and railroad access. I estimate the treatment model by multivariate lasso logistic regression (Friedman et al., 2010).

The squared loss in Eq. (5) is weighted by estimated propensity scores, \hat{w}_{it} , to place more emphasis on the loss for the values in \mathcal{O} most similar to the counterfactual values in terms of pre-treatment outcomes and covariates. Consistent estimation of the ATT (1) relies on the correct specification of the outcome model (3) under the assumption of exogeneity (4). It does not rely on the correct specification of the treatment model (6); although, the propensity scores from estimating this model are intended to help balance treated and control units in terms of the pre-treatment outcomes and covariates when fitting the matrix completion model on the factual data.

The algorithm for solving Eq. (5) iteratively replaces missing values with those recovered from a singular value decomposition of the matrix (Mazumder et al., 2010). Once \mathbf{L} , γ , and δ have been estimated, we can predict the counterfactual values for the treated units in the post-treatment period by

$$\hat{Y}(\infty)_{it} = \hat{L}_{it} + \hat{\gamma}_i + \hat{\delta}_t, \quad \forall (i, t) \notin \mathcal{O}.$$

3.4. Simulation studies

In simulation studies described in Section A2 in the online Appendix, I conduct two sets of simulation studies to assess the performance of the matrix completion estimator: the first on generated data in which we control the ground-truth treatment effects; and the second on empirical data, where the focus is on time periods where no treatment effects are expected. The comparison estimators are evaluated on their ability to recover the ground-truth ATT averaged over the counterfactual period, $\tau_{\infty a'_i}$ (2). The comparison estimators include two versions of the matrix completion estimator: with (MC-W) and without (MC) a treatment propensity-weighted loss function. The DID estimator is a regression of outcomes on treatment and unit and time fixed effects. The SCM is a regression of the pre-treatment

outcomes of each treated unit on the control unit outcomes during the same periods, with the restrictions of no intercept and non-negative regression weights that sum to one. The SCM with lasso (SCM-L1) relaxes the zero-intercept and weight restrictions and estimates the counterfactual outcomes for each treated unit by lasso regression. I provide the exact form of these estimators in Section A3.

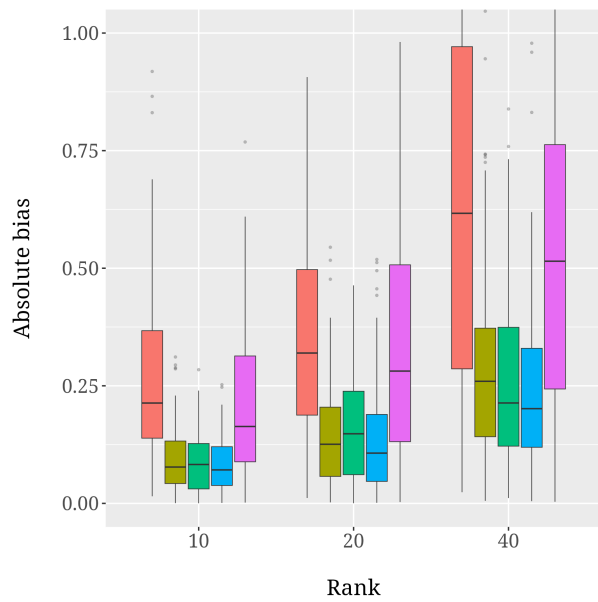
3.4.1. Generated data

Figure 2 provides box and whisker plots summarizing the median, the first and third quartiles, and outlying points of the distribution of the absolute bias — i.e., absolute difference between $\hat{\tau}_{\infty a'_i}$ and the actual $\tau_{\infty a'_i}$ — and the variance of 399 block bootstrap replicates of $\hat{\tau}_{\infty a'_i}$ for the first set of simulations on generated data. The absolute bias and bootstrap variance increase across all estimators as the rank of L_{it} increases, which underscores the importance of the low-rank assumption. The matrix completion estimators and the SCM estimator yield the lowest absolute bias and bootstrap variance, regardless of rank, whereas the DID and SCM-L1 estimators struggle with higher bias and variance. When the rank is high relative to N and T , MC-W exhibits lower absolute bias and bootstrap variance relative to the unweighted MC estimator, suggesting that propensity score weighting may mitigate the impact of increased model complexity on estimator bias and variance.

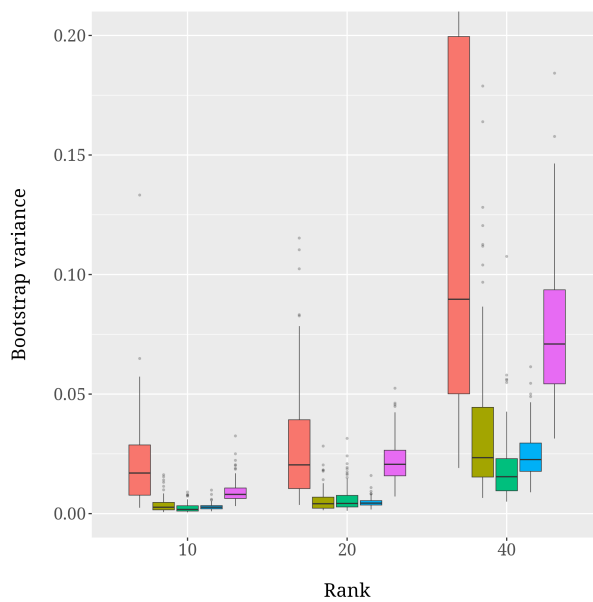
3.4.2. Empirical data

In the second set of simulations focusing on the empirical data, I leverage the fact that the true treatment effect is null in the pre-treatment period. I first discard the post-treatment data, and for each of 1000 simulation runs, randomly select half of the control units to be treated and impute their missing values following a placebo treatment time randomly chosen from $\{a'_i, \dots, T\}$, varying a'_i .

In simulation studies on the state government finance datasets, described in Section 4,



(A) Absolute bias



(B) Bootstrap variance

Figure 2: Absolute bias and bootstrap variance in generated data, varying the rank of L_{it} .
Estimator: DID; MC; MC-W; SCM; SCM-L1.

the matrix completion estimators generally maintain lower absolute bias than the DID estimator while exhibiting higher bias when compared to SCM estimators, across all placebo a'_i ratios for both expenditure and revenue datasets (Figure A2a). In terms of bootstrap variance, matrix completion estimators demonstrate results on par with synthetic control estimators and markedly outperform the DID estimator (Figure A2b). There is not much efficiency gain from propensity-weighting the matrix completion estimator in the empirical data simulation studies because, unlike the generated data simulations, treatment is assigned at random, rather than as a function of covariates.

In each of three datasets common to the synthetic control literature, the matrix completion estimators outperform DID and the SCM estimators by minimizing absolute bias (Figure A3) and bootstrap variance (Figure A4) across the different ratios of the placebo a'_i to T . Together, the simulation results support the preferential use of the MC-W estimator in the application.

4. Application: Homestead policy and state size

In order to estimate causal impacts of homestead policies on state size, I create measures of total expenditure and revenue collected from the records of 48 state governments during the period of 1789 to 1932 (Sylla et al., 1993), 16 state governments during the period of 1933 to 1937 (Sylla et al., 1995a,b), and U.S. Census special reports for the period of 1902 to 2008, covering 48 states (Haines, 2010; U.S. Census Bureau, 2010). The expenditure measure includes state government spending on education, social welfare programs, and transportation. The revenue measure incorporates state government income streams such as tax revenue and non-tax revenue such as land sales.

The expenditure and revenue data pre-processing steps are as follows. Removing years with zero or near-zero variance results in outcome matrices consisting of $T = 203$ observations for $N = 48$ states, 30 of which are treated. The outcome data are inflation-adjusted

according to the U.S. Consumer Price Index (Williamson, 2017) and scaled by the total free adult male population in the decennial census (Haines, 2010). I impute the outcome values that are missing due to lack of data collection using multiple imputation by chained equations (MICE, Buuren and Groothuis-Oudshoorn, 2010). Figure A5 visualizes the extent of the missing data in the entire dataset by state and treatment group, where 40% of values in the dataset are missing (29.9% and 10.1% missing in the control and treated groups, respectively). The majority of the outcome data for treated states prior to the treatment time were missing and have been imputed. To address the concern that the choice of imputation method can influence the estimated treatment effects, Table A1 evaluates the sensitivity of the causal estimates to alternative imputation methods. Lastly, I log-transform the data to alleviate exponential effects.

The staggered treatment implementation setting is appropriate in this application because a_i varies across states that were exposed to homesteads following the passage of the HSA. I aggregate to the state level approximately 1.46 million individual land patent records authorized under the HSA. Using these records, which are made available by the BLM (General Land Office, 2017), I determine that the earliest homestead entries occurred in 1869 in about half of the western frontier states, about seven years following the enactment of the HSA. In 1872, the first homesteads were filed in southern PLS. Figure A6 shows how each state is categorized in the empirical analysis, as a PLS (treated group) or state-land state (control group), as well as the year of the earliest initial homestead entry for the PLS, which informs staggered treatment implementation.

I include the following covariates in the conditioning set of the treatment model (6): per-capita spending or revenue prior to 1869; the ratio of slaves to the total population in 1860; and the ratio of free African-Americans, Native Americans, or Whites to the total non-slave population in 1860; average farm sizes in 1860 and average farm values in 1850 and 1860 (Haines, 2010); and the state-level share of total miles of operational railroad track per square mile, which I calculate by overlaying the railroad track map over historical county

borders (Atack, 2013). These pre-treatment covariates control for selective migration to more agriculturally productive land, and for differences in the accessibility and availability of frontier lands. Bustos (2017, p. 45) finds that the prevalence of slavery in 1860 is an important predictor of available homestead lands, and reasons that the covariate acts as a proxy for the presence of large plantations.

4.1. Accounting for bias

Potential sources of bias include violations of the assumptions of exogeneity (4), no interference, no-anticipation, or invariance to treatment history. The exogeneity assumption would be violated if the error term in the outcome model (3) is correlated with the initial treatment time. While this assumption is not directly testable, the no-treatment evaluation on pre-treatment data reported in Section 5.2 provide indirect evidence that the exogeneity assumption is not violated. Additionally, the simulation results on the state government finances datasets reported in Section 3.4 demonstrate that propensity-weighting the loss function improves the consistency of the matrix completion estimator.

A second potential source of bias arises from interference, or the assumption that control units are unaffected by the effects of treatment. While the no interference assumption cannot directly be tested, it is likely in the present application that the outcomes of state-land states were indirectly affected by the out-migration of homesteaders from frontier states. When assuming the absence of interference, the use of indirectly affected states as control units would underestimate treatment effects because it would make the counterfactual and factual treated unit observations in the post-treatment period more similar. Interference might also arise if state-land state governments increase public investments in order to dissuade workers from migrating to the frontier in the first place. The historical evidence, however, suggests that labor-scarce frontier states were more strongly motivated to attract migrants and stimulate population growth than long-settled state-land states (Engerman and Sokoloff, 2005). For example, the adoption of compulsory primary education laws and

support for public education in general in western states has been considered as a means to attract potential migrants to the frontier (Meyer et al., 1979; Bandiera et al., 2018). Interference arising from competition among state governments would also underestimate the effect of treatment.

A third potential source of bias arises in violations of the no-anticipation or invariance to treatment history assumptions. The no-anticipation assumption would be violated if there were anticipatory effects on the size of frontier state governments prior to the initial homestead entries. Anticipatory effects are plausible since the first homestead entries occurred in 1869 in western PLS, six years after the HSA went into effect. In Section 5.2, I conduct a no-treatment evaluation on the pre-treatment data and vary the placebo initial treatment year. The estimated placebo ATT is nonsignificant for most settings, which is direct evidence of no anticipatory effects. The invariance to treatment history assumption rules out variation in treatment effects by the initial treatment time, but does not rule out causal effects of treatment duration. In Section 5.1, I explore whether causal effects on state size differ with respect to year of initial homestead entry.

Lastly, bias may result from misspecification of imputation models. The imputation procedure assumes that after controlling for the available state government finances data, the missing values are Missing At Random (MAR). There are reasons to believe the data are not MAR, which could result in biased estimates. For example, the timing of a state's admission to the Union, which affects the extent of its missing values, may be determined by unobserved political and demographic variables rather than meeting a population threshold. It is impossible to distinguish whether data are MAR or missing based on unobserved variables, given the observed data (Sterne et al., 2009). In Section 5, the sensitivity of causal estimates to two alternative imputation methods is evaluated, indicating that the choice of imputation method alters the conclusions in one out of the two scenarios examined, as detailed in Table A1.

5. Matrix completion estimates

I estimate the causal impacts of the initial treatment year on the state government finances of the treated units (i.e., PLS). Specifically, I fit the MC-W estimator on the entirety of factual outcomes to recover the counterfactual outcomes of the treated units had they not been exposed to treatment. The top panel of Figure 3 compares the average log per-capita state government expenditure of treated units and control units along with the predicted average expenditure of treated units. The dashed vertical line represents the year of the earliest homestead entry, $a'_i = \min_{1 \leq i \leq N_T} a_i = 1869$. The treated unit outcomes are generally higher than those of the control units in the pre-treatment period, whereas there is little difference between the treated and control unit outcomes, on average, in the post-treatment period.

The difference between the observed and predicted treated unit outcomes, which is the quantity $\hat{\tau}_{t, \infty a'_i}$ described in Eq. (1), corresponds to the estimated per-period ATT. These per-period average causal impacts are plotted in the bottom panel, with 95% normal interval confidence intervals estimated by calculating the standard error of the distribution of block bootstrap replicates of $\hat{\tau}_{t, \infty a'_i}$. The bootstrap replicates are constructed by block resampling the columns (i.e, time dimension) of the observed outcomes, in order to preserve temporal dependence structure of the original data (Davison and Hinkley, 1997; Politis and White, 2004), and obtaining a set of point estimates from 999 resamples. The estimated per-period effects for both outcomes are essentially zero during the pre-treatment period and within the bounds of the bootstrap confidence intervals, which demonstrates that the model is closely fitting the pre-treatment period observations. By 1876, after most PLS had been exposed to homesteads, homestead exposure decreases per-capita state government expenditure by about 0.13 log points, and the trajectory of estimated causal impacts remains negative for the rest of the time-series, although the confidence intervals for the per-period effects all contain zero. Similar patterns are observed when the outcome is

per-capita state government revenue (Figure A7).

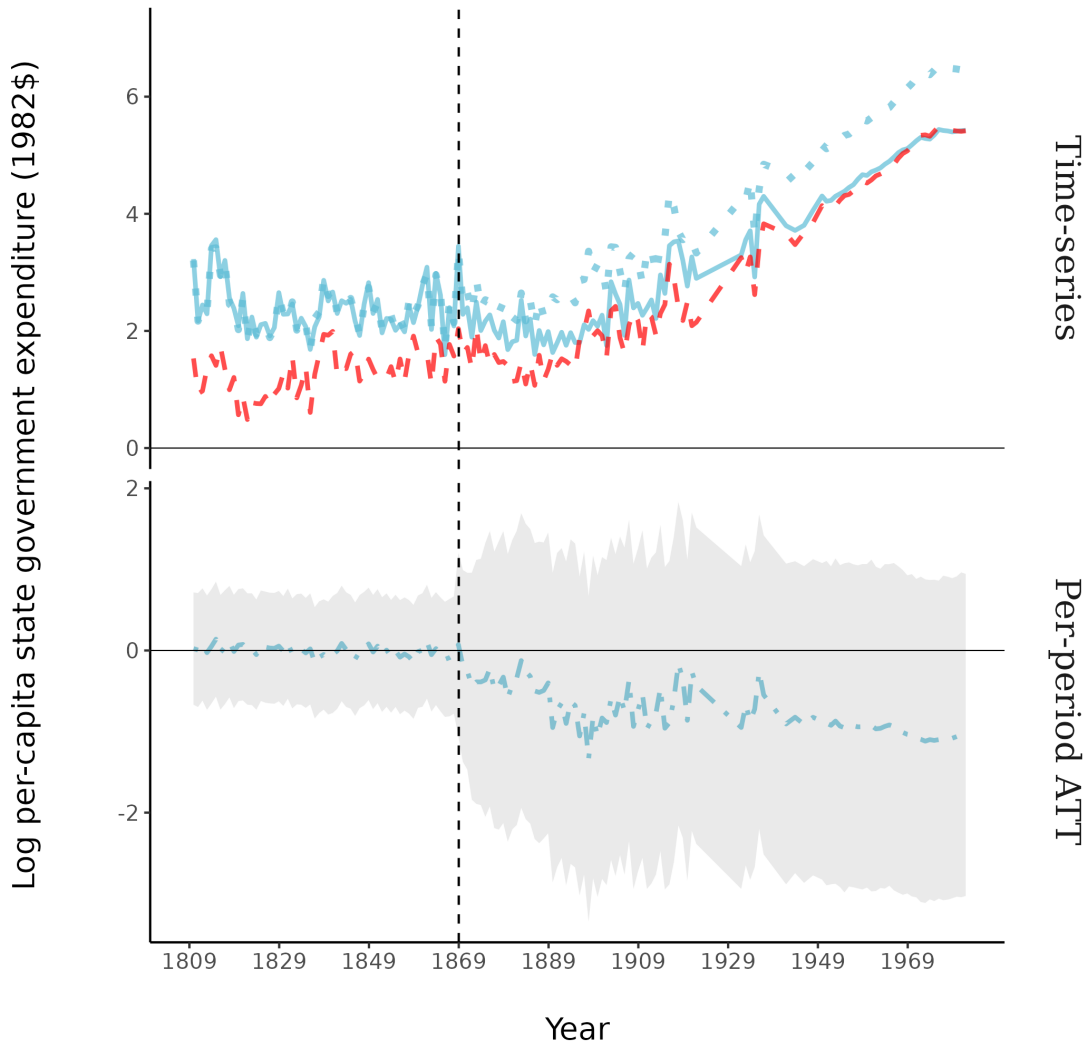


Figure 3: Matrix completion estimates of the effect of the year of initial homestead entry (1869; dashed vertical line) on state government expenditure, 1809 to 1982: —, factual treated; - - - , factual control; ·····, counterfactual treated; - · - · , $\hat{\tau}_{t, \infty a'_i}$.

To infer the overall effect of treatment, I estimate the ATT averaged over the counterfactual period of 1869 to 2008, and report the point estimates and bootstrap standard errors in the second and third columns of Table 1. The estimates reveal that per-capita state government expenditure would have been 0.84 [0.61, 1.08] log points lower had the PLS never been exposed to homesteads. Relative to the observed log per-capita state ex-

penditure of the PLS in the same period, the point estimate represents a decrease of 0.22%. The estimated ATT on per-capita state government revenue is similar.

I compare the MC-W estimates with DID and SCM estimates, also reported in the second and third columns of Table 1. Figure A8 illustrates the parallel trends assumption in DID analysis by depicting the log per-capita state government expenditures and revenues over time for both treated and control groups, showing similar trajectories up to the treatment year of 1869, which supports the validity of the DID approach. The point estimates from the binary DID estimator are slightly larger and within the confidence intervals of the MC-W estimates. The ATT estimates from the SCM estimator are positive, but not statistically significant.

Table A1 presents counterfactual period estimates on differently imputed datasets. When estimated on data with missing outcome values imputed by MICE with classification and regression trees (CART) as the imputation method, rather than predictive mean matching (the default method), the conclusions drawn from the estimates do not change. However, when estimating on data with missing values imputed by an expectation-maximization (EM) algorithm based method, the MC-W estimates are much larger in magnitude and no longer statistically significant. In interpreting the results presented in Table A1, it is important to reflect on the implications of imputing a majority of the outcomes in treated states prior to the treatment period. The differences in ATT estimates, especially under the EM imputation method, underscore the sensitivity of our results to the imputation of missing data.

5.1. Treatment effect heterogeneity

Recall that under staggered treatment implementation, the time of initial treatment a_i varies across states that were exposed to homesteads, and that the year of the earliest homestead entry among the treated units a'_i is used to calculate the ATT. Also recall that the SHA opened land for homesteading in the South under the same stipulations as the

Table 1: Estimates of the ATT averaged over the counterfactual period of 1869 to 2008 and bootstrap standard errors (in parentheses).

	<i>All PLS</i>		<i>Southern PLS</i>		<i>Western PLS</i>	
	Expenditure	Revenue	Expenditure	Revenue	Expenditure	Revenue
DID	-0.90 (0.32)	-0.95 (0.27)	-0.69 (0.39)	-0.55 (0.35)	-0.94 (0.45)	-1.04 (0.51)
MC-W	-0.84 (0.24)	-0.79 (0.25)	-0.63 (0.30)	-0.42 (0.30)	-0.88 (0.37)	-0.86 (0.39)
SCM	0.13 (0.12)	0.17 (0.14)	-0.09 (0.21)	-0.14 (0.23)	0.17 (0.26)	0.23 (0.27)
SCM-L1	0.13 (0.17)	0.17 (0.19)	-0.08 (0.23)	-0.11 (0.21)	0.17 (0.27)	0.22 (0.29)

HSA, which opened land for homesteading in the western frontier. The results above set a'_i at 1869, which is the earliest homestead entry that occurred in the western PLS. Among southern PLS, the earliest homestead entry occurred in 1872. Thus, there is a substantive interest in determining whether there is a differential effect of the year of initial homestead entry on state size based on region. Conducting a sub-group analysis by region also allows us to detect potential violations in the assumption of invariance to treatment history, since most of the western PLS are treated for a longer period than the southern PLS.

The last four columns of Table 1 decomposes the counterfactual period estimates by calculating the ATT with respect to the region of the PLS. The MC-W estimates show that the main effect on all of the PLS (second and third columns) is driven mainly by the effect on the Western PLS. The estimated effect size on the southern PLS are comparatively smaller in magnitude and significant for the effect on state government expenditure, -0.63 [-0.94, -0.33], but not on revenue, -0.42 [-0.73, -0.12]. These results provide indirect evidence that the assumption of invariance to treatment history is not violated since the conclusions drawn from the main estimates are generally unchanged.

5.2. No-treatment evaluation

To assess whether the estimated effects are attributable to the year of initial homestead entry rather than other policy changes or spurious errors during the same period, I conduct

a no-treatment evaluation by discarding the post-treatment period observations from the state government finances data and re-running the analysis on the pre-treatment data, when no treatment effect is expected (i.e., the ATT is zero).

Table 2 reports placebo ATT estimates and block bootstrap standard errors on each outcome, considering $t = \{1, \dots, a'_i - \Delta\}$ as the pre-treatment period, with $\Delta \in \{1, 10, 25\}$. Across all estimators, the standard error decreases with a larger Δ , reflecting the uncertainty of estimating causal effects in shorter (placebo) post-treatment periods. Compared to the binary DID and SCM estimators, MC-W exhibits lower standard errors in all settings and lower bias in three of the six settings. The placebo ATT estimates from the MC-W estimator is significant only when the outcome is state government expenditure and Δ is 10 or 25. Similar patterns are observed when conducting no-treatment evaluations on differently imputed datasets (Table A2). These placebo results bolster the usage of the MC-W estimator in the application, and provide evidence supporting the plausibility of the exogeneity and no-anticipation assumptions.

Table 2: Placebo ATT estimates and bootstrap standard errors (in parentheses).

	Expenditure			Revenue		
	$\Delta = 1$	$\Delta = 10$	$\Delta = 25$	$\Delta = 1$	$\Delta = 10$	$\Delta = 25$
DID	-0.61 (0.43)	-0.38 (0.24)	-0.39 (0.22)	-0.63 (0.48)	-0.70 (0.42)	-0.67 (0.34)
MC-W	-0.26 (0.32)	-0.44 (0.16)	-0.47 (0.15)	-0.33 (0.38)	-0.26 (0.21)	-0.28 (0.16)
SCM	0.94 (0.45)	0.86 (0.24)	0.76 (0.19)	0.66 (0.47)	0.52 (0.28)	0.64 (0.22)
SCM-L1	0.48 (0.41)	0.34 (0.22)	0.48 (0.16)	0.61 (0.39)	0.39 (0.26)	0.20 (0.21)

6. Continuous DID estimation

The matrix completion approach estimates the impact of a binary exposure to treatment on a continuous outcome. In the application, however, a continuous form of treatment is available in the number of homestead patents. The model below is a continuous version of the DID estimator described in Section A3.1, where the first difference comes from

variation in the date of initial exposure to homesteads, and the second difference comes from variation in the intensity of homestead entries:

$$Y_{it} = \xi_i + \psi_t + \zeta W_{it} + \phi(W_{it} \cdot H_{it}) + \beta X_{ip} + v_{it}. \quad (7)$$

The model includes state and year fixed effects, $\{\xi_i\}_{i=1}^N$ and $\{\psi_t\}_{t=1}^T$, respectively. The covariate X_{ip} controls for average farm sizes and values, and railroad access, when the outcome Y_{it} is log per-capita state government spending or revenue; X_{ip} controls for average farm values when Y_{it} is land inequality. The continuous treatment variable H_{it} measures the log of the per-capita number of patents issued under the HSA in state i and year t . The coefficient of interest corresponds to the interaction term, ϕ , which represents the average causal effect of exposure to homesteads. A least squares estimator for ϕ is given by

$$\arg \min_{\phi, \xi_i, \psi_t, \zeta, \beta} \sum_{i=1}^N \sum_{t=1}^T (Y_{it} - \xi_i - \psi_t - \zeta W_{it} - \phi(W_{it} \cdot H_{it}) - \beta X_{ip})^2. \quad (8)$$

6.1. Estimates on state size and land inequality

Table 3 reports DID estimates of the average causal effect of exposure to log per-capita homestead patents, with 95% confidence intervals constructed using 999 state-stratified bootstrap samples. The estimates indicate that a 1% increase in log per-capita homesteads decreases log per-capita state government spending or revenue by about 4%. The point estimates are considerably smaller in magnitude — albeit, in the same direction — as the per-period MC-W estimates presented in Section 5. The bootstrap confidence intervals around the DID estimates are considerably more narrow than those for the MC-W estimates displayed in the bottom panels of Figures 1 and A7, and are potentially overoptimistic due to serial correlation in the DID regression errors (Bertrand et al., 2004). The estimates on state size are insensitive to the method used for imputing expenditure or revenue values that are missing due to nonresponse (Table A3).

Table 3: Continuous DID estimates of the effect of (log) per-capita homestead patents on state size or land inequality.

	Expenditure	Revenue	Land inequality
Treatment effect ($\hat{\phi}$)	-0.04 (0.002)	-0.04 (0.002)	-0.0007 (0.0003)
Adjusted R ²	0.538	0.540	0.801
N	8,618	8,618	463

The third column of Table 3 presents DID estimates of the impact of log per-capita homesteads on land inequality at the state-level during the period of 1870 to 1950. Since land inequality is measured every decennial, I aggregate homesteads to the next decennial year; e.g., the number of homesteads measured in 1880 is the total for the years 1871 to 1880. Average farm values are included in the regression as a proxy for agricultural productivity, which might be associated with farm sizes approaching ideal scale and therefore land inequality. I estimate that homesteads significantly decreased land inequality in frontier states: a 1% increase in log per-capita homesteads lowers state-level land inequality by about 10^{-5} points.

The direction of the point estimate is consistent with the study of Bustos (2017, p. 3), who conducts a county-level DID analysis and shows treatment based on terciles of Homestead Act acres reduced land inequality measured by the Gini coefficient over a similar post-treatment period, although the magnitude of the estimate in the present work is substantially smaller. The comparatively small coefficient implies that homestead policies did not fundamentally alter the long-run distribution of landownership, which may be explained by qualitative evidence that suggests homestead policies were exploited by land speculators and natural resource companies and that the rents from public land were appropriated by the private sector.

7. Conclusion

The findings of this paper signify that mid-nineteenth century homestead policies had long-lasting impacts that can potentially explain contemporary differences in state government size. Estimates using matrix completion with a binary treatment and DID with continuous treatment evidence that homestead policies had significant and negative impacts on state government expenditure and revenue that lasted a century following their implementation. This finding is in line with recent work documenting the adverse impact of homestead policies on the economic development of regions exposed to homesteading.

I explore land inequality as a possible causal mechanism underlying the relationship between homestead policies and state size, which is closely related to state capacity. First, I provide evidence of a positive relationship between land inequality and state government finances and that the slope of correlation increases at higher levels of inequality. Nonlinearities in the relationship between inequality and state capacity can arise in theoretical models that incorporate economic differences in political influence: greater income inequality reduces government spending and investments in state capacity when elites have a monopoly on political power, however when inequality gets too high, the poor can impose redistribution through majority voting. Second, I present continuous DID estimates that reveal per-capita homesteads significantly lowered land inequality in frontier states; although, the magnitude of the effect is negligible. This finding is in line with previous empirical work showing that exposure to homesteads decreased land inequality. The failure to fundamentally alter the long-run distribution of landownership may be explained by qualitative evidence that suggests homestead policies were *de facto* corporate welfarism often exploited by land speculators and corporations to amass land and resources during early capitalist expansion.

This paper makes a methodological contribution by extending the matrix completion method for causal effect estimation in panel data with staggered treatment adoption to

allow for propensity score weighting of the loss function. The matrix completion estimator with propensity score weighting outperforms regression-based estimators such as the synthetic control method and difference-in-differences in simulation studies and a no-treatment evaluation. This methodological contribution holds implications for policy evaluation, offering a more accurate tool for understanding the effects of policies over time and place.

Funding details

This work was supported by the National Science Foundation under Grants DGE 1106400 and TG-SES 180010.

Supplemental online material

The online Appendix includes simulation results, and describes model specifications and implementation details for each of the comparison estimators used in the simulations. It includes descriptive figures on the extent of missing data in the state government finances datasets, and reports the results of sensitivity analyses on differently imputed datasets. It also includes figures for matrix completion estimates of treatment exposure on state government revenue, a diagnostic plot for the DID parallel trends assumption, and bivariate regression estimates of the relationship between land inequality and state government finances.

Data and code

Data and R code to reproduce the results of the paper are available at <https://github.com/jvpoulos/homesteads>.

References

- Abadie, A., Diamond, A., and 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, 493–505.
- Acemoglu, D. and Robinson, J. A. (2008), “Persistence of Power, Elites, and Institutions,” *American Economic Review*, 98, 267–293.
- Agarwal, A., Shah, D., Shen, D., and Song, D. (2021), “On Robustness of Principal Component Regression,” *arXiv:1902.10920*.
- Allen, D. W. (1991), “Homesteading and Property Rights; Or, “How the West Was Really Won”,” *The Journal of Law and Economics*, 34, 1–23.
- Amjad, M., Shah, D., and Shen, D. (2018), “Robust synthetic control,” *The Journal of Machine Learning Research*, 19, 802–852.
- Ansell, B. and Samuels, D. J. (2015), *Inequality and Democratization: An Elite Competition Approach*, Cambridge: Cambridge University Press.
- Arkhangelsky, D., Athey, S., Hirshberg, D. A., Imbens, G. W., and Wager, S. (2021), “Synthetic Difference in Differences,” *arXiv:1812.09970*.
- Atack, J. (2013), “On the Use of Geographic Information Systems in Economic History: The American Transportation Revolution Revisited,” *The Journal of Economic History*, 73, 313–338.
- Athey, S., Bayati, M., Doudchenko, N., Imbens, G., and Khosravi, K. (2021), “Matrix completion methods for causal panel data models,” *Journal of the American Statistical Association*, 1–15.
- Athey, S. and Imbens, G. W. (2021), “Design-based analysis in difference-in-differences settings with staggered adoption,” *Journal of Econometrics*.
- Bai, J. and Ng, S. (2021), “Matrix Completion, Counterfactuals, and Factor Analysis of Missing Data,” *arXiv:1910.06677*.
- Bandiera, O., Mohnen, M., Rasul, I., and Viarengo, M. (2018), “Nation-building through compulsory schooling during the age of mass migration,” *The Economic Journal*, 129, 62–109.
- Bazzi, S., Fiszbein, M., and Gebresilashe, M. (2020), “Frontier culture: The roots and persistence of ‘rugged individualism’ in the United States,” *Econometrica*, 88, 2329–2368.
- Ben-Michael, E., Feller, A., and Rothstein, J. (2018), “The Augmented Synthetic Control Method,” *arXiv:1811.04170*.

- Ben-Michael, E., Feller, A., and Rothstein, J. (2019), “Synthetic Controls with Staggered Adoption,” *arXiv:1912.03290*.
- Benabou, R. (2000), “Unequal Societies: Income Distribution and the Social Contract,” *American Economic Review*, 96–129.
- Bensel, R. F. (1990), *Yankee Leviathan: the Origins of Central State Authority in America, 1859-1877*, Cambridge: Cambridge University Press.
- Bertrand, M., Duflo, E., and Mullainathan, S. (2004), “How Much Should We Trust Differences-in-Differences Estimates?” *The Quarterly Journal of Economics*, 119, 249–275.
- Besley, T. and Persson, T. (2009), “The Origins of State Capacity: Property Rights, Taxation and Politics,” *American Economic Review*, 99, 1218–1244.
- (2010), “State Capacity, Conflict, and Development,” *Econometrica*, 78, 1–34.
- Boix, C. (2003), *Democracy and Redistribution*, Cambridge: Cambridge University Press.
- Bustos, N. A. L. (2017), “Essays on the effects of the Homestead Act on Land Inequality and Human Capital, the effects of Land Redistribution on Crop Choice, and the effects of Earthquakes on Birth Outcomes,” Ph.D. thesis, University of Warwick, available at <https://core.ac.uk/download/pdf/186333176.pdf>.
- Buuren, S. v. and Groothuis-Oudshoorn, K. (2010), “MICE: Multivariate imputation by chained equations in R,” *Journal of Statistical Software*, 10, 1–68.
- Callaway, B. and Sant’Anna, P. H. (2020), “Difference-in-Differences with multiple time periods,” *Journal of Econometrics*.
- Danhof, C. H. (1941), “Farm-making costs and the “safety valve”: 1850-60,” *Journal of Political Economy*, 49, 317–359.
- Davison, A. C. and Hinkley, D. V. (1997), *Bootstrap Methods and their Application*, vol. 1, Cambridge University Press.
- Deverell, W. F. (1988), “To Loosen the Safety Valve: Eastern Workers and Western Lands,” *The Western Historical Quarterly*, 19, 269–285.
- Doudchenko, N. and Imbens, G. W. (2016), “Balancing, Regression, Difference-In-Differences and Synthetic Control Methods: A Synthesis,” *arXiv:1610.07748*.
- Engerman, S. L. and Sokoloff, K. L. (2005), “The Evolution of Suffrage Institutions in the New World,” *The Journal of Economic History*, 65, 891–921.
- Fan, J., Masini, R. P., and Medeiros, M. C. (2021), “Do We Exploit all Information for Counterfactual Analysis? Benefits of Factor Models and Idiosyncratic Correction,” *arXiv:2011.03996*.

- Ferrie, J. P. (1997), “Migration to the Frontier in Mid-Nineteenth Century America: A Re-Examination of Turner’s ‘Safety Valve’,” Available at <https://faculty.wcas.northwestern.edu/~fe2r/papers/munich.pdf>.
- Friedman, J., Hastie, T., and Tibshirani, R. (2010), “Regularization paths for generalized linear models via coordinate descent,” *Journal of Statistical Software*, 33, 1.
- Frymer, P. (2014), “‘A Rush and a Push and the Land Is Ours’: Territorial Expansion, Land Policy, and U.S. State Formation,” *Perspectives on Politics*, 12, 119.
- (2017), *Building an American Empire: The Era of Territorial and Political Expansion*, Princeton, NJ: Princeton University Press.
- Galor, O., Moav, O., and Vollrath, D. (2009), “Inequality in Landownership, the Emergence of Human-Capital Promoting Institutions, and the Great Divergence,” *The Review of Economic Studies*, 76, 143–179.
- García-Jimeno, C. and Robinson, J. A. (2008), “The myth of the frontier,” in *Understanding Long-Run Economic Growth: Geography, Institutions, and the Knowledge Economy*, Chicago, IL: University of Chicago Press, pp. 49–88.
- Gates, P. W. (1936), “The Homestead Law in an Incongruous Land System,” *The American Historical Review*, 41, 652–681.
- (1940), “Federal Land Policy in the South 1866-1888,” *The Journal of Southern History*, 6, 303–330.
- (1942), “The Role of the Land Speculator in Western Development,” *The Pennsylvania Magazine of History and Biography*, 66, 314–333.
- (1979), “Federal Land Policies in the Southern Public Land States,” *Agricultural History*, 53, 206–227.
- Gates, P. W. and Bogue, A. G. (1996), *The Jeffersonian dream: Studies in the history of American land policy and development*, Albuquerque, NM: University of New Mexico Press.
- General Land Office (2017), *General Land Office Records Automation*, Bureau of Land Management, Washington, DC, available at: <https://glorerecords.blm.gov/>.
- Gobillon, L. and Magnac, T. (2016), “Regional policy evaluation: Interactive fixed effects and synthetic controls,” *Review of Economics and Statistics*, 98, 535–551.
- Goodman-Bacon, A. (2021), “Difference-in-differences with variation in treatment timing,” *Journal of Econometrics*.
- Haines, M. R. (2010), “Historical, Demographic, Economic, and Social Data: The United States, 1790-2002,” Ann Arbor, MI: Inter-university Consortium for Political and Social Research [distributor], 2010-05-21. doi.org/10.3886/ICPSR02896.v3.

- Hoffnagle, W. (1970), “The Southern Homestead Act: Its Origins and Operation,” *The Historian*, 32, 612–629.
- Lanza, M. L. (1999), *Agrarianism and Reconstruction Politics: The Southern Homestead Act*, Baton Rouge, LA: LSU Press.
- Long, J. H. (1995), “Atlas of Historical County Boundaries,” *The Journal of American History*, 81, 1859–1863.
- Mattheis, R. and Raz, I. T. (2021), “There’s no such thing as free land: the Homestead Act and economic development,” Available at https://extranet.sioe.org/uploads/sioe2021/mattheis_raz.pdf.
- Mazumder, R., Hastie, T., and Tibshirani, R. (2010), “Spectral regularization algorithms for learning large incomplete matrices,” *Journal of Machine Learning Research*, 11, 2287–2322.
- Meltzer, A. H. and Richard, S. F. (1981), “A Rational Theory of the Size of Government,” *Journal of Political Economy*, 89, 914–927.
- Meyer, J. W., Tyack, D., Nagel, J., and Gordon, A. (1979), “Public education as nation-building in America: Enrollments and bureaucratization in the American states, 1870–1930,” *American Journal of Sociology*, 85, 591–613.
- Murtazashvili, I. (2013), *The Political Economy of the American Frontier*, Cambridge: Cambridge University Press.
- Neyman, J. (1923), “On the application of probability theory to agricultural experiments,” *Annals of Agricultural Sciences*, 51, reprinted in Splawa-Neyman et al. (1990).
- Perotti, R. (1993), “Political Equilibrium, Income Distribution, and Growth,” *The Review of Economic Studies*, 60, 755–776.
- Politis, D. N. and White, H. (2004), “Automatic Block-Length Selection for the Dependent Bootstrap,” *Econometric Reviews*, 23, 53–70.
- Poulos, J. and Zeng, S. (2021), “RNN-based counterfactual prediction, with an application to homestead policy and public schooling,” *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 70, 1124–1139.
- Poulos, J. et al. (2023), “Gender Gaps in Frontier Entrepreneurship? Evidence from 1901 Oklahoma Land Lottery Winners,” *Journal of Historical Political Economy*, 2, 611–634.
- Ramcharan, R. (2010), “Inequality and Redistribution: Evidence from U.S. Counties and States, 1890–1930,” *The Review of Economics and Statistics*, 92, 729–744.
- Rubin, D. B. (1990), “Comment: Neyman (1923) and Causal Inference in Experiments and Observational Studies,” *Statistical Science*, 5, 472–480.

- Saint-Paul, G. and Verdier, T. (1993), “Education, Democracy and Growth,” *Journal of Development Economics*, 42, 399–407.
- Shanks, T. R. (2005), “The Homestead Act: A Major Asset-Building Policy in American History,” *Inclusion in the American Dream: Assets, Poverty, and Public Policy*, 20–41.
- Splawa-Neyman, J., Dabrowska, D., Speed, T., et al. (1990), “On the application of probability theory to agricultural experiments. Essay on principles. Section 9,” *Statistical Science*, 5, 465–472.
- Sterne, J. A., White, I. R., Carlin, J. B., Spratt, M., Royston, P., Kenward, M. G., Wood, A. M., and Carpenter, J. R. (2009), “Multiple imputation for missing data in epidemiological and clinical research: Potential and pitfalls,” *BMJ*, 338.
- Sylla, R. E., Legler, J. B., and Wallis, J. (1993), “Sources and Uses of Funds in State and Local Governments, 1790-1915: [United States],” Ann Arbor, MI: Inter-university Consortium for Political and Social Research, 2017-05-21. doi.org/10.3886/ICPSR06304.v1.
- (1995a), “State and Local Government [United States]: Sources and Uses of Funds, Census Statistics, Twentieth Century [Through 1982],” Ann Arbor, MI: Inter-university Consortium for Political and Social Research, 2017-05-21. doi.org/10.3886/ICPSR06304.v1.
- (1995b), “State and Local Government [United States]: Sources and Uses of Funds, State Financial Statistics, 1933-1937,” Ann Arbor, MI: Inter-university Consortium for Political and Social Research, 2017-05-21. http://doi.org/10.3886/ICPSR06306.v1.
- Turner, F. J. (1956), *The Significance of the Frontier in American History*, Ithaca, NY: Cornell University Press.
- U.S. Census Bureau (2010), “Data Base on Historical Finances of Federal, State and Local Governments,” Available at: <https://www.census.gov/programs-surveys/gov-finances/data/historical-data.html>.
- Vollrath, D. (2013), “Inequality and School Funding in the Rural United States, 1890,” *Explorations in Economic History*, 50, 267–284.
- Williamson, S. H. (2017), “Seven ways to compute the relative value of a U.S. dollar amount, 1774 to present,” Available at <http://MeasuringWorth.com>.
- Xiong, R. and Pelger, M. (2020), “Large Dimensional Latent Factor Modeling with Missing Observations and Applications to Causal Inference,” *arXiv:1910.08273*.
- Xu, Y. (2017), “Generalized Synthetic Control Method: Causal Inference with Interactive Fixed Effects Models,” *Political Analysis*, 25, 57–76.
- Ziblatt, D. (2008), “Does Landholding Inequality Block Democratization? A Test of the “Bread and Democracy” Thesis and the Case of Prussia,” *World Politics*, 60, 610–641.

Online Appendix for State-Building through Public Land Disposal? An Application of Matrix Completion for Counterfactual Prediction

Jason Poulos

December 29, 2023

Abstract

The online Appendix includes maps describing the treatment status of states (Figures 1 and A6); bivariate regression estimates of the relationship between land inequality and state government finances (Figure A1); a write-up of the design and results of the simulation studies (Section A2); and the model specifications and implementation details for each of the comparison estimators (Section A3). Lastly, Section A4 includes tables and figures for the extent of missing data in the state government finances datasets (Figure A5); matrix completion estimates of treatment exposure on state government revenue (Figure A7); a diagnostic plot for the DID parallel trends assumption (Figure A8); and reports the results of sensitivity analyses on differently imputed datasets (Table A1).

A1. Empirical application: Descriptive statistics

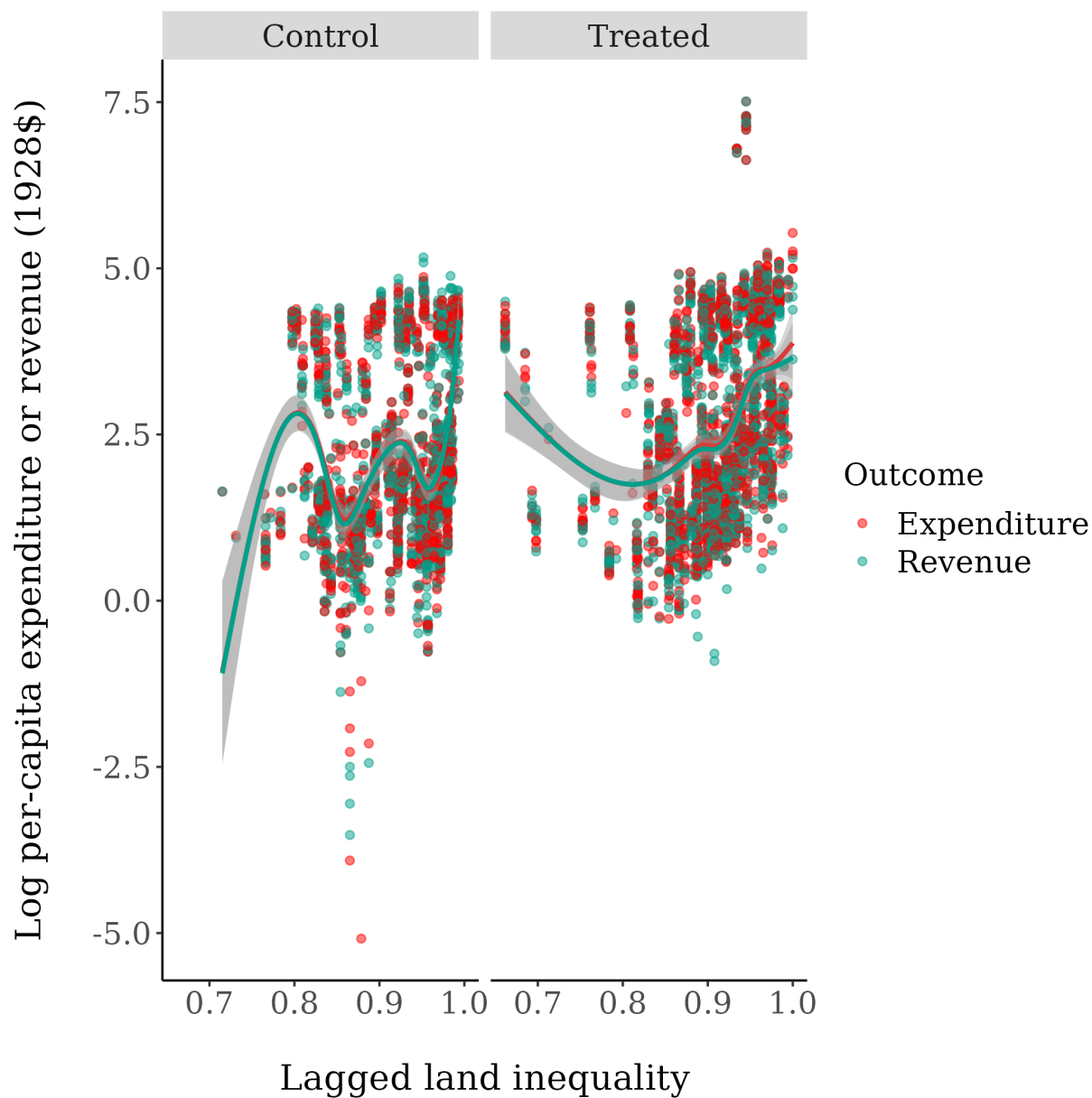


Figure A1: Land inequality (lagged by 10 years) vs. log per-capita state government revenue and expenditure, 1860-1950. Each point is a state-year observation. Lines represent generalized additive model fits to the data for the two outcomes and shaded regions represent corresponding 95% confidence intervals. The model is fit separately on control states (i.e., state-land states) and treated states (i.e., PLS).

A2. Simulations

A2.1. Generated data

In the first set of simulations, I generate potential outcomes under control according to the outcome model (3). For each of 100 trial runs, I generate the low-rank matrix L_{it} as the product of factors V_{tr} and factor loadings U_{ir} . The factors are drawn independently from a multivariate normal distribution with means and variances of 1, and covariances of 0.2; the factor loadings are simulated from a first-order autoregressive model with slope coefficient $\phi = 0.3$. I generate ground-truth propensity scores according to the treatment model (6), where the simulated covariate is also generated using a first-order autoregressive model with $\phi = 0.3$. I focus on the case of a square matrices, $N \times T = 60 \times 60$. For each trial run, I sample $N_T = 30$ treated units, and use the ground-truth propensity scores as sampling weights. Each treated unit is assigned an initial treatment period, with the first treatment time set to $a'_i = 3$.

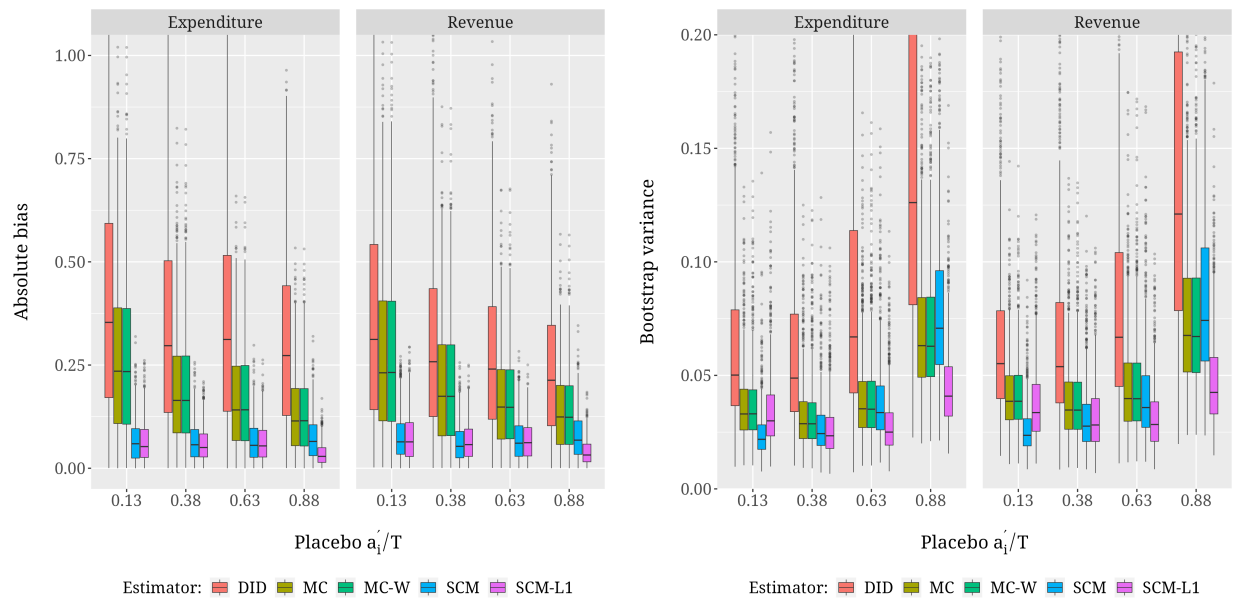
A2.2. Empirical data

A2.2.1. State government finances datasets

I evaluate the performance of the matrix completion estimator on the state government expenditure and revenue datasets described in Section 4 of the main paper, discarding the treated units and using the control units ($N = 18$, $T = 203$). Figures A2a and A2b present box and whisker plots summarizing the absolute bias and variance, respectively, across 1000 simulation runs on the state government finances datasets. The x axis in each figure is the ratio of the placebo initial treatment time to the number of periods in the placebo data, so higher values represent more training data.

A2.2.2. Synthetic control datasets

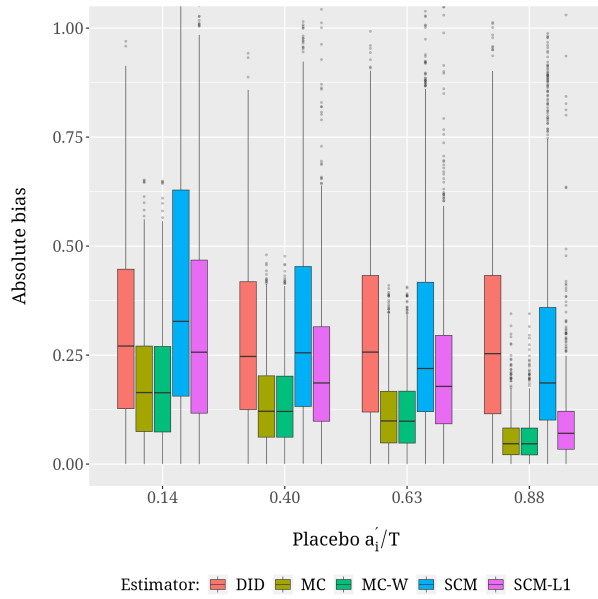
Next, I evaluate the performance of the matrix completion estimator on three datasets common to the synthetic control literature, with the actual treated unit removed from each dataset. The three synthetic control datasets originate from Abadie and Gardeazabal's [2003] study of the economic impact of terrorism in the Basque Country during the late 1960s ($N = 16$, $T = 43$); Abadie et al.'s [2010] study of the effects of a large-scale tobacco control program implemented in California in 1988 ($N = 38$, $T = 31$); and Abadie et al.'s [2015] study of the economic impact of the 1990 German reunification on West Germany ($N = 16$, $T = 44$). Figure A3 provides box and whisker plots of the absolute bias for each estimator on the three synthetic control datasets.



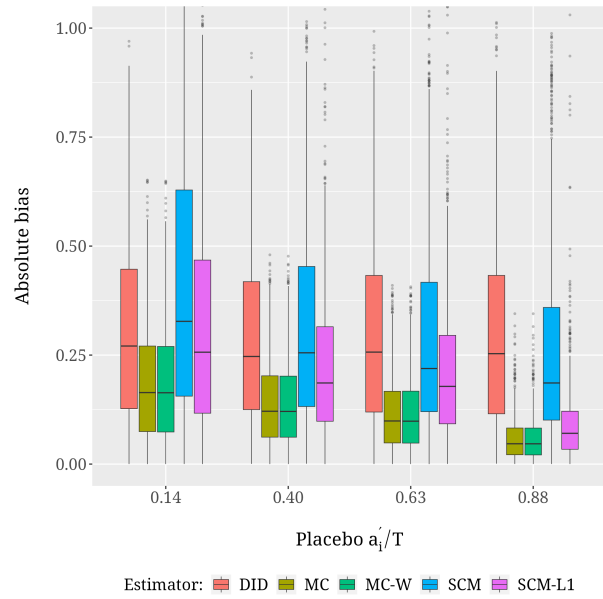
(A) Absolute bias

(B) Bootstrap variance

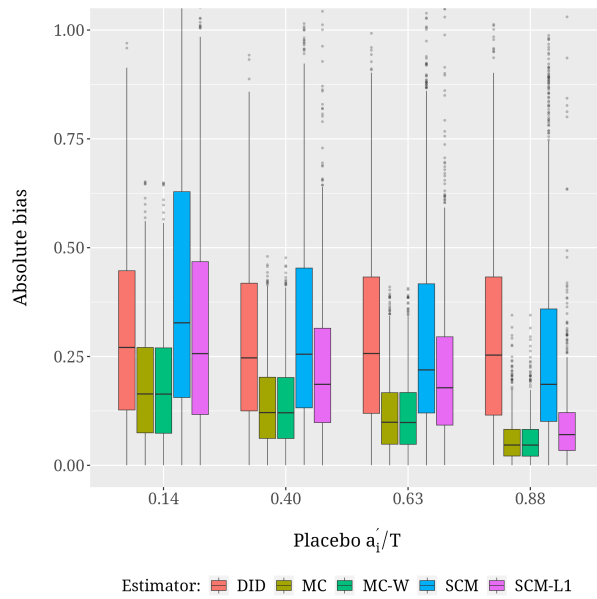
Figure A2: Absolute bias and bootstrap variance for state government finances datasets, varying the placebo a'_i .



(A) Basque Country terrorism

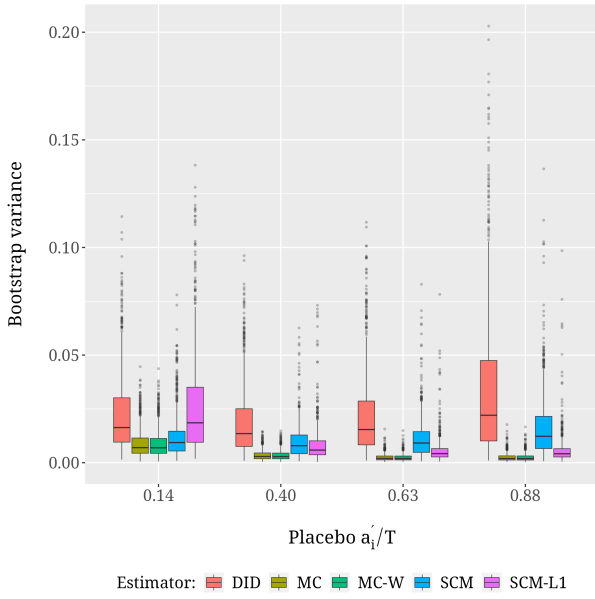


(B) California smoking ban

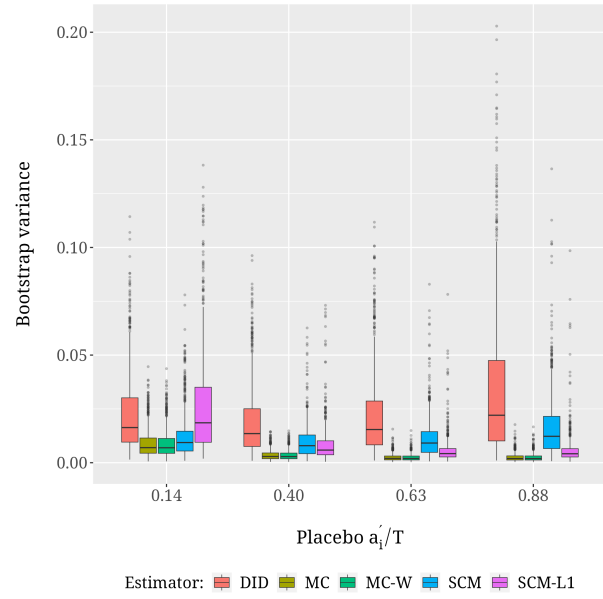


(C) West German reunification

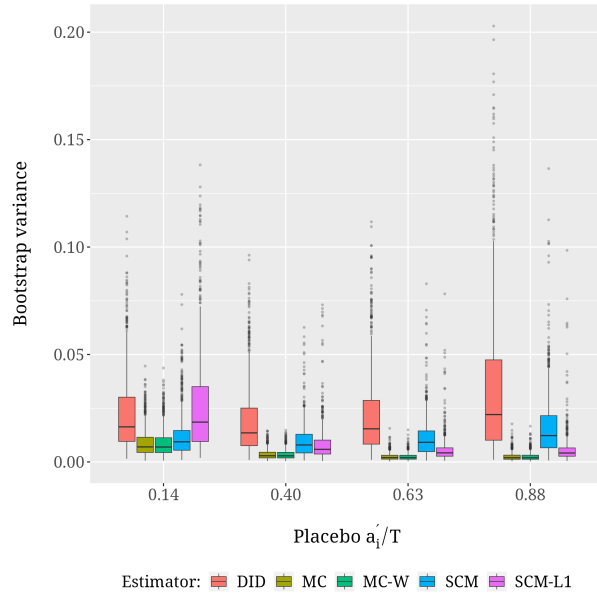
Figure A3: Absolute bias for synthetic control datasets, varying the placebo a'_i .



(A) Basque Country terrorism



(B) California smoking ban



(C) West German reunification

Figure A4: Bootstrap variance for synthetic control datasets, varying the placebo a'_i .

A3. Benchmark estimators

The following estimators are used for comparison in the no treatment evaluation (Section 2) and main estimates (Section 5).

A3.1. Difference-in-differences (DID)

The DID model with binary treatment is specified by Athey and Imbens (2021). The outcome under control is modeled as:

$$Y(\infty)_{it} = \xi_i + \psi_t + \tau W_{it} + v_{it}, \quad (1)$$

where $\{\xi_i\}_{i=1}^N$ are unit-specific fixed effects and $\{\psi_t\}_{t=1}^T$ are time-specific fixed effects. The model is estimated by least squares:

$$\arg \min_{\tau, \xi_i, \psi_t} \sum_{i=1}^N \sum_{t=1}^T (Y_{it} - \xi_i - \psi_t - \tau W_{it})^2. \quad (2)$$

I use the implementation of DID in the `MCPanel` R package (Athey et al., 2017).

A3.2. Matrix completion

Matrix completion methods attempt to impute missing values by solving a convex optimization problem via nuclear norm minimization, even when relatively few values are observed in the data matrix (Candès and Recht, 2009; Candès and Plan, 2010; Mazumder et al., 2010; Recht, 2011). I implement two versions of matrix completion estimated by via nuclear norm regularized least squares (Athey et al., 2021): with (MC-W) and without (MC) an propensity-weighted loss function. The MC-W outcome model and estimation equation is specified in Eqs. (3) and (5) of the main paper, respectively. The MC outcome model is the same Eq. (3) of the main paper, while the estimation equation does not weight the loss function:

$$\arg \min_{\mathbf{L}, \gamma, \delta} \left[\frac{1}{|\mathcal{O}|} \sum_{(i,t) \in \mathcal{O}} \left(Y_{it} - L_{it} - \gamma_i - \delta_t \right)^2 + \lambda_L \|\mathbf{L}\|_{\star} \right]. \quad (3)$$

Both versions are implemented using an extended version of the `MCPanel` package that includes an option for propensity-weighting the loss function.¹

A3.3. Synthetic control method (SCM)

The synthetic control method (SCM, Abadie et al., 2010) compares a single treated unit with a synthetic control that combines the outcomes of multiple control units on the basis of their pre-intervention similarity with the treated unit.

¹Available at <https://github.com/jvpoulos/MCPanel>.

Doudchenko and Imbens (2016) and Athey et al. (2021) show that the SCM can be interpreted as regressing the pre-treatment outcomes of a single treated unit on the control unit outcomes during the same periods. The parameters estimated on the controls are then used to predict the counterfactual outcomes for a single treated unit, $i = 0$:

$$\hat{Y}_{0t} = \sum_{i=1}^N \hat{\omega}_i Y_{it}, \quad \forall t = \alpha', \dots, T,$$

$$\text{where } \hat{\omega} = \arg \min_{\omega} \sum_{s=1}^{\alpha'} \left(Y_{0s} - \sum_{i=1}^N \omega_i Y_{0s} \right)^2, \quad \text{s.t. } w_i \geq 0, \sum w_i = 1. \quad (4)$$

A separate model is subsequently fit for each i, \dots, N_t treated units. Note that Eq. (4) imposes the restrictions of the original SCM, namely zero intercept and non-negative regression weights that sum to one. I use the `MCPanel` implementation with default settings, except I bound gradient values within $[-5, 5]$ in order to facilitate convergence.

A3.4. SCM with lasso (SCM-L1)

The SCM with lasso (SCM-L1, Doudchenko and Imbens, 2016; Athey et al., 2021) relaxes the zero-intercept and weight restrictions of Eq. (4). The counterfactual outcomes for treated unit $i = 0$ is:

$$\hat{Y}_{0t} = \hat{\mu} + \sum_{i=1}^N \hat{\omega}_i Y_{it}, \quad \forall t = \alpha', \dots, T,$$

$$\text{where } (\hat{\mu}, \hat{\omega}) = \arg \min_{\mu, \omega} \sum_{s=1}^{\alpha'} \left(Y_{0s} - \mu - \sum_{i=1}^N \omega_i Y_{0s} \right)^2, \quad (5)$$

where a separate model is fit for each i, \dots, N_t treated units. Intuitively, the generalized SCM is a convex combination of control units with intercept μ and weight ω_i for control units i, \dots, N . The model is fit with $N + 1$ predictors, including the number of control units and the intercept, and α', \dots, T observations.

Eq. (5) is estimated by lasso regression (Tibshirani, 1996; Tibshirani et al., 2012) in order to reduce the relative size of the predictor set. I use the `MCPanel` implementation with the strength of the lasso penalty selected among 30 possible values by 5-fold cross-validation.

A4. Empirical application: tables & figures

- Figure A6 shows how each state is categorized in the empirical analysis, and the year of the earliest initial homestead entry for the treated states.
- Figure A5 visualizes the extent of the missing values in the state government finances datasets.
- Figure A1 plots bivariate regression estimates of the relationship between land inequality and state government finances.
- Figure A7 plots matrix completion estimates of treatment exposure on state government revenue.
- Figure A8 visualizes the average outcomes of treated and control groups for the purpose of assessing the parallel trends assumption for DID estimation with a binary treatment variable.
- In the no treatment evaluation and in the main analysis, the missing values are imputed by multivariate imputation by chained equations (MICE, Azur et al., 2011) with predictive mean matching as the imputation method, implemented using the `mice` R package (Buuren and Groothuis-Oudshoorn, 2010). Table A1 reports the results of sensitivity analyses on differently imputed datasets using the following two alternative imputation methods:

MICE-CART MICE with classification and regression trees (CART) from the `rpart` package (Therneau and Atkinson, 2022) as the imputation method;

EM An expectation–maximization (EM) algorithm based method for imputing missing values in multivariate normal time series, implemented using the `mtsdi` package with default settings (Therneau and Atkinson, 2022).

In order to avoid train-test set contamination, each imputation method is fit only on the outcome values of the control units. Table A2 reports the estimated ATT on differently imputed placebo datasets, and Table A3 reports treatment effect estimates using the continuous DID model (Section 6) on differently imputed datasets.

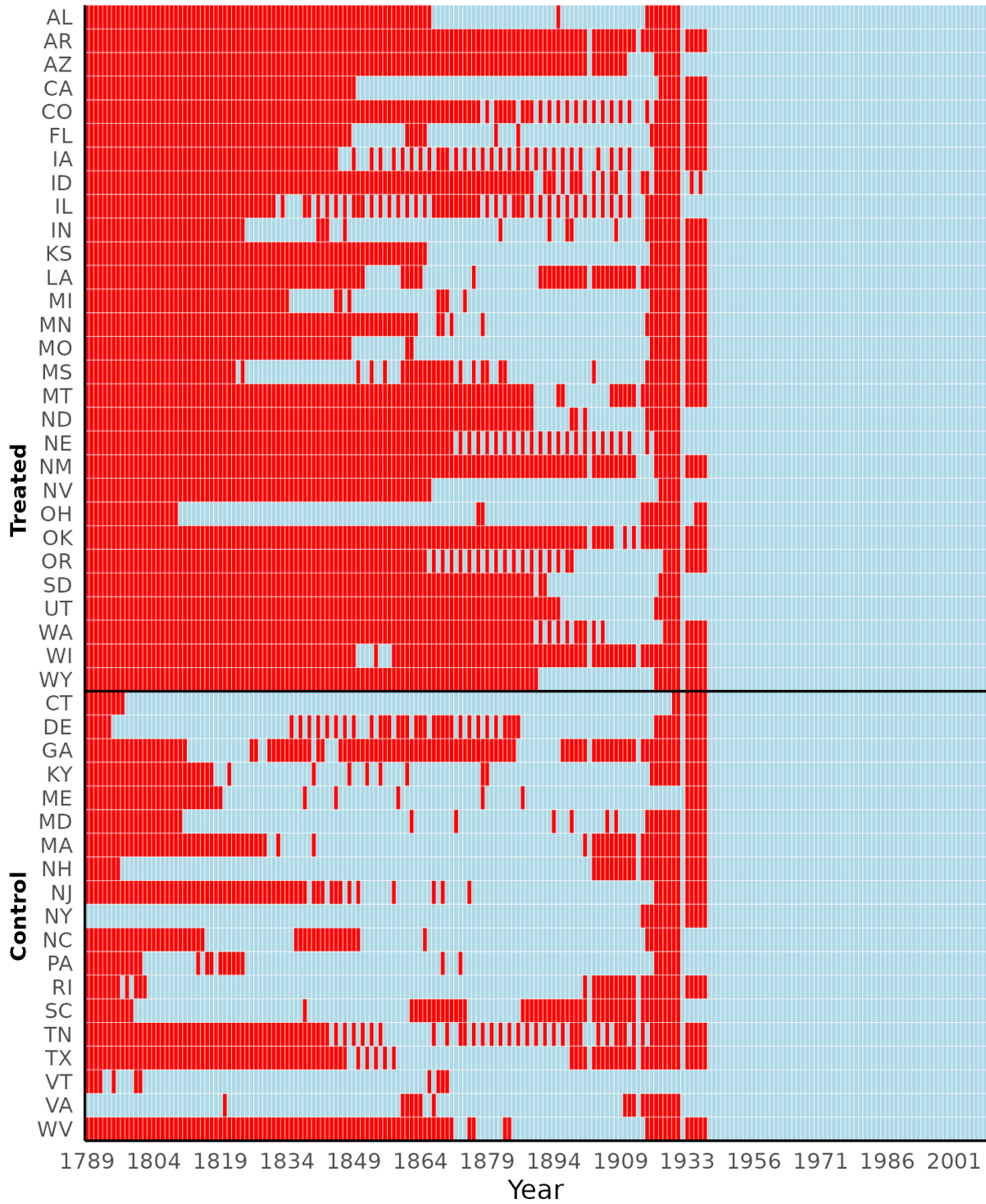


Figure A5: State-year observations in the state government finances datasets that are missing (■) or present (□).

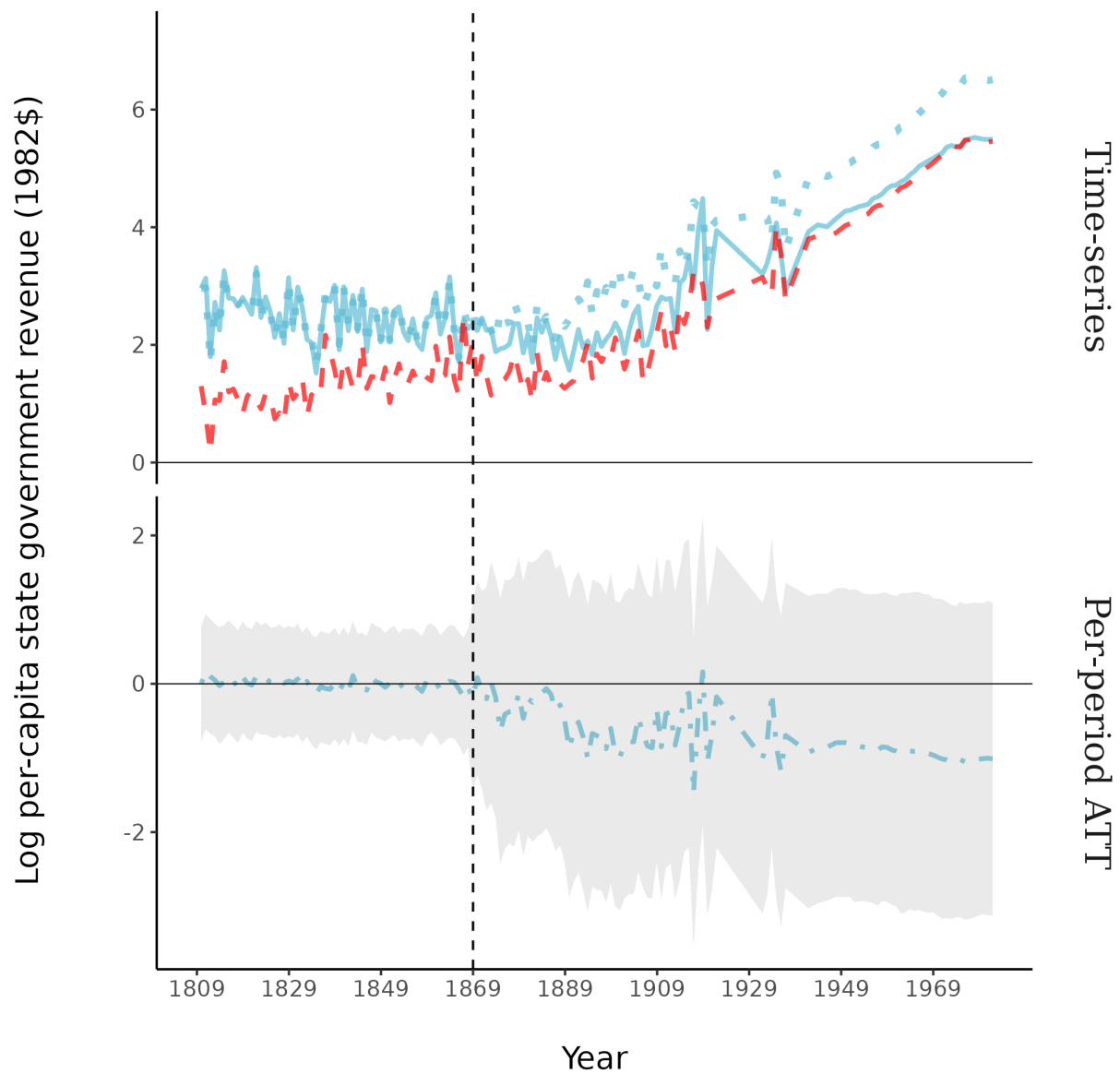


Figure A7: Matrix completion estimates of the effect of the year of initial homestead entry (1869; dashed vertical line) on state government revenue, 1809 to 1982: —, factual treated; - - -, factual control; ·····, counterfactual treated; — · —, $\hat{\tau}_{t, \infty a'_i}$.

Table A1: Estimates of the ATT averaged over the counterfactual period (2) and bootstrap standard errors (in parentheses) on differently imputed datasets.

	<i>All PLS</i>		<i>Southern PLS</i>		<i>Western PLS</i>	
	Expenditure	Revenue	Expenditure	Revenue	Expenditure	Revenue
<i>Imputation method: MICE-CART</i>						
DID	-0.66 (0.26)	-0.73 (0.27)	-0.53 (0.37)	-0.76 (0.41)	-0.69 (0.39)	-0.73 (0.39)
MC-W	-0.69 (0.21)	-0.63 (0.21)	-0.52 (0.32)	-0.66 (0.33)	-0.72 (0.33)	-0.62 (0.32)
SCM	0.13 (0.14)	0.08 (0.14)	-0.11 (0.22)	-0.16 (0.24)	0.18 (0.25)	0.13 (0.18)
SCM-L1	0.16 (0.16)	0.11 (0.15)	-0.11 (0.22)	-0.15 (0.22)	0.19 (0.25)	0.16 (0.24)
<i>Imputation method: EM</i>						
DID	-2.29 (0.82)	-1.48 (0.55)	-2.11 (0.85)	-1.45 (0.87)	-2.32 (0.96)	-1.48 (0.93)
MC-W	-2.04 (0.60)	-1.34 (0.46)	-1.79 (0.71)	-1.26 (0.74)	-2.09 (0.80)	-1.36 (0.79)
SCM	0.32 (0.38)	0.11 (0.31)	0.12 (0.44)	-0.33 (0.54)	0.36 (0.50)	0.20 (0.36)
SCM-L1	0.21 (0.42)	0.20 (0.33)	-0.01 (0.48)	-0.21 (0.53)	0.26 (0.54)	0.26 (0.59)

Table A2: Placebo ATT estimates and bootstrap standard errors (in parentheses) on differently imputed datasets.

	Expenditure			Revenue		
	$\Delta = 1$	$\Delta = 10$	$\Delta = 25$	$\Delta = 1$	$\Delta = 10$	$\Delta = 25$
<i>Imputation method: MICE-CART</i>						
DID	-0.35 (0.51)	-0.60 (0.39)	-0.51 (0.28)	-0.04 (0.53)	-0.61 (0.29)	-0.49 (0.25)
MC-W	-0.19 (0.43)	-0.52 (0.23)	-0.45 (0.17)	-0.08 (0.38)	-0.52 (0.19)	-0.27 (0.15)
SCM	0.51 (0.53)	0.16 (0.24)	0.79 (0.52)	0.10 (0.43)	-0.21 (0.24)	0.09 (0.17)
SCM-L1	0.47 (0.48)	-0.15 (0.22)	0.12 (0.18)	0.84 (0.44)	0.05 (0.20)	0.22 (0.17)
<i>Imputation method: EM</i>						
DID	-0.81 (0.87)	-1.27 (0.80)	-1.11 (0.63)	-1.75 (1.81)	-2.11 (1.42)	-1.76 (1.05)
MC-W	-0.41 (0.54)	-0.68 (0.32)	-0.64 (0.30)	-1.71 (1.35)	-0.59 (0.81)	-1.26 (0.64)
SCM	0.89 (0.90)	0.65 (0.66)	0.79 (0.52)	-0.85 (1.47)	-0.16 (1.18)	0.38 (0.88)
SCM-L1	-0.19 (0.55)	-0.16 (0.39)	0.26 (0.36)	-1.18 (0.88)	-0.04 (0.72)	-0.59 (0.65)

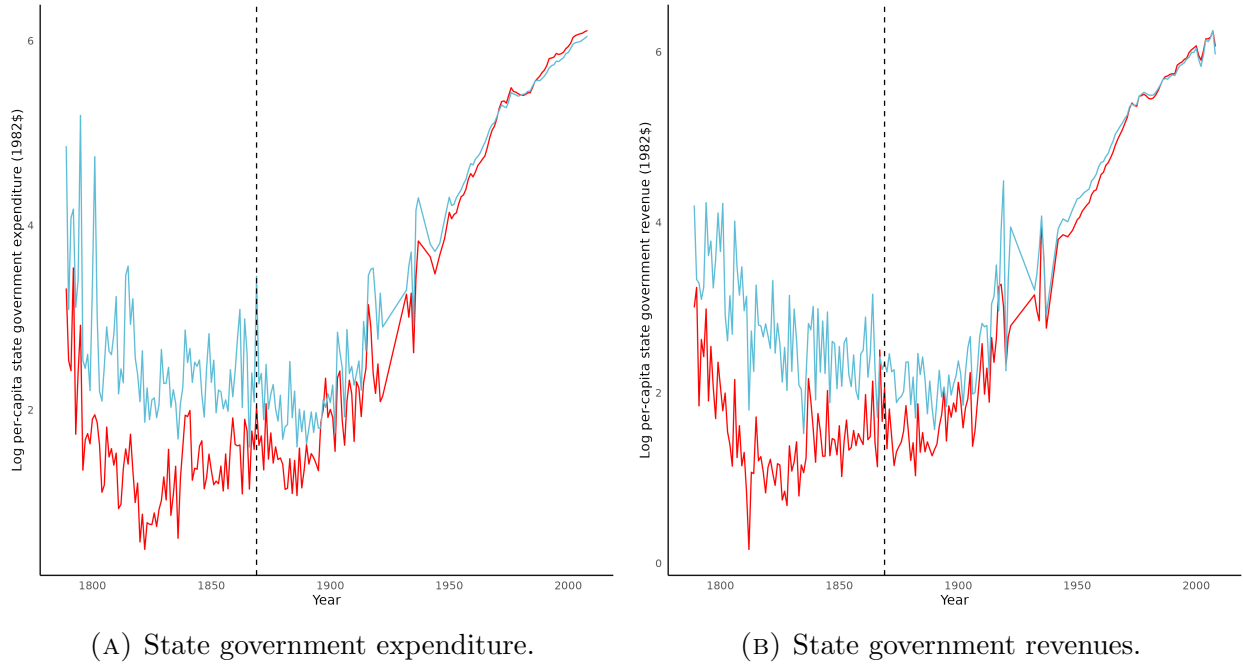


Figure A8: Visually assessing the parallel trends assumption for DID estimation with a binary treatment variable: —, factual treated; —, factual control. The dashed vertical line represents the earliest treatment year, 1869.

Table A3: Continuous DID estimates of the effect of (log) per-capita homestead patents on differently imputed state government finance datasets.

<i>Imputation method:</i>	MICE-CART		EM	
<i>Outcome:</i>	Expenditure	Revenue	Expenditure	Revenue
Treatment effect ($\hat{\phi}$)	-0.04 (0.002)	-0.04 (0.002)	-0.04 (0.002)	-0.06 (0.003)
Adjusted R ²	0.568	0.577	0.452	0.320
N	8,618	8,618	8,618	8,618

References

- Abadie, A., Diamond, A., and 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, 493–505.
- (2015), “Comparative Politics and the Synthetic Control Method,” *American Journal of Political Science*, 59, 495–510.
- Abadie, A. and Gardeazabal, J. (2003), “The Economic Costs of Conflict: A Case Study of the Basque Country,” *The American Economic Review*, 93, 113–132.
- Athey, S., Bayati, M., Doudchenko, N., Imbens, G., and Khosravi, K. (2017), “MCPanel: Matrix Completion Methods for Causal Panel Data Models,” Available at: <https://github.com/susanathey/MCPanel>.
- (2021), “Matrix completion methods for causal panel data models,” *Journal of the American Statistical Association*, 1–15.
- Athey, S. and Imbens, G. W. (2021), “Design-based analysis in difference-in-differences settings with staggered adoption,” *Journal of Econometrics*.
- Azur, M. J., Stuart, E. A., Frangakis, C., and Leaf, P. J. (2011), “Multiple imputation by chained equations: What is it and how does it work?” *International Journal of Methods in Psychiatric Research*, 20, 40–49.
- Buuren, S. v. and Groothuis-Oudshoorn, K. (2010), “MICE: Multivariate imputation by chained equations in R,” *Journal of Statistical Software*, 10, 1–68.
- Candes, E. J. and Plan, Y. (2010), “Matrix completion with noise,” *Proceedings of the IEEE*, 98, 925–936.
- Candès, E. J. and Recht, B. (2009), “Exact matrix completion via convex optimization,” *Foundations of Computational Mathematics*, 9, 717.
- Doudchenko, N. and Imbens, G. W. (2016), “Balancing, Regression, Difference-In-Differences and Synthetic Control Methods: A Synthesis,” *arXiv:1610.07748*.
- Mazumder, R., Hastie, T., and Tibshirani, R. (2010), “Spectral regularization algorithms for learning large incomplete matrices,” *Journal of Machine Learning Research*, 11, 2287–2322.
- Recht, B. (2011), “A simpler approach to matrix completion,” *Journal of Machine Learning Research*, 12, 3413–3430.
- Therneau, T. M. and Atkinson, E. J. (2022), “An introduction to recursive partitioning using the RPART routines,” Available at: <https://stat.ethz.ch/R-manual/R-patched/library/rpart/doc/longintro.pdf>.

Tibshirani, R. (1996), “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 267–288.

Tibshirani, R., Bien, J., Friedman, J., Hastie, T., Simon, N., Taylor, J., and Tibshirani, R. J. (2012), “Strong rules for discarding predictors in lasso-type problems,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 74, 245–266.