

# NONNEGATIVE TENSOR COMPLETION FOR DYNAMIC COUNTERFACTUAL PREDICTION ON COVID-19 PANDEMIC

BY YAOMING ZHEN<sup>1,a</sup>  AND JUNHUI WANG<sup>2,b</sup> 

<sup>1</sup>*School of Data Science, City University of Hong Kong, [ayzhen8-c@my.cityu.edu.hk](mailto:ayzhen8-c@my.cityu.edu.hk)*

<sup>2</sup>*Department of Statistics, The Chinese University of Hong Kong, [junhuiwang@cuhk.edu.hk](mailto:junhuiwang@cuhk.edu.hk)*

The COVID-19 pandemic has been a worldwide health crisis for the past three years, casting unprecedented challenges for policymakers in different countries and regions. While one country or region can only implement one social mobility restriction policy at a given time, it is of great interest for policy makers to decide whether to elevate or delevate the restriction policy from time to time. This article proposes a novel nonnegative tensor completion method to predict the potential counterfactual outcomes of multifaceted social mobility restriction policies over time. The proposed method builds upon a low-rank tensor decomposition of the pandemic data, which also explicitly characterizes the ordinal nature of the mobility restriction strength and the smooth trend of the pandemic evolution over time. Its application to the COVID-19 pandemic data reveals some interesting facts regarding the impact of social mobility restriction policy on the spread of the virus. The effectiveness of the proposed method is also supported by its asymptotic estimation consistency and extensive numerical experiments on the synthetic datasets.

**1. Introduction.** The Coronavirus Disease 2019 (COVID-19) pandemic is caused by a novel coronavirus identified in 2019, called the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2). According to the COVID-19 Dashboard at <https://coronavirus.jhu.edu/map.html> by the Center for Systems Science and Engineering at Johns Hopkins University, the pandemic has quickly spread over 200 countries and regions, with about 435 million confirmed cases and 5.9 million deaths as of March 1, 2022. Social distancing has been proved to be an effective measure to prevent the virus from spreading (Yan et al. (2021); Sinha and Chakraborty (2022); Quick, Dey and Lin (2021); Ye et al. (2021)), yet at the great cost of sacrificing the economy. Many governments have implemented various social mobility restriction policies to enforce social distancing, to a certain extent, so as to reduce the spread of the COVID-19 pandemic. It is thus of great significance for administration units or policymakers to understand the potential outcomes, or counterfactuals, of multifaceted public policies at a particular pandemic stage so that suitable policies can be adopted for the confronting situations.

In literature a great deal of existing works have been developed to investigate the counterfactual causal effects (Höfler (2005); Pearl (2009)) of different treatments or interventions. Imputation by regression is one of the most popular methods (Johansson, Shalit and Sontag (2016); Agarwal, Shah and Shen (2022); Fan, Masini and Medeiros (2022)), which requires a fully observed preintervention period and invariant relationships among the treated individuals before and after different interventions. More recently, Athey et al. (2021), Bai and Ng (2021), and Poulos et al. (2022) formulate the counterfactual prediction as a matrix completion problem, which requires that the missing entries must locate in a single block of the data matrix after appropriately permuting the rows and columns. This is equivalent to assuming

---

Received December 2022; revised April 2023.

*Key words and phrases.* Causal inference, imputation, informative missing, latent factor, tensor decomposition.

the existence of fully observed rows and columns. Another line of research related to counterfactual prediction is to estimate the average treatment effect (Athey, Imbens and Wager (2018); Fogarty (2020); Yadlowsky et al. (2021)) and derive the corresponding individualized treatment rules (Qi et al. (2020); Mo, Qi and Liu (2021); Mo and Liu (2022)).

However, the aforementioned methods are generally not suitable for the COVID-19 pandemic data due to a number of reasons. First, as one country or region only implements one social mobility restriction policy for a given time period, it leads to only one observation for each country or region at a given time. Second, the similarity of pandemic evolution patterns between two countries may vary over time due to different policy trajectories. Third, public policies are often multifaceted and multilevel, while existing imputation methods are restricted to binary intervention. For example, the number of quarantine days and the number of people allowed for social gathering are constantly adjusted under different pandemic situations.

In this paper we aim to conduct the counterfactual prediction of daily mortality rates of different countries or regions under various social mobility restriction policies. We formulate the task as a nonnegative tensor completion problem (Chen et al. (2019); Zhang and Ng (2022)) and propose to estimate the data tensor with a nonnegative tensor possessing multiple important features. Specifically, the estimated tensor admits a low-rank CANDECOMP/PARAFAC (CP) decomposition, where all the factor matrices of the expected data tensor to be nonnegative, the factor matrix in the policy mode is ordinal, and that in the time mode possesses a fusion structure. The nonnegativity requirement is imposed to ensure the expected data tensor to be nonnegative, corresponding to the nonnegative daily mortality rates. The ordinal structure is of great significance to characterize the ordinality of the mobility restriction strength of different public policies, and the fusion structure plays a crucial role in modeling the smooth trend of the pandemic evolution over time. The nonnegative tensor completion problem is then formulated in a least square format and subsequently optimized via an efficient alternative gradient descent algorithm.

The main contribution of this paper is three-fold. First, we develop a dynamic nonnegative tensor completion framework to predict the daily counterfactual mortality rates caused by the COVID-19 pandemic under various potential intervention policies in different countries or regions. The counterfactual prediction leads to a number of interesting results. For instance, it provides an accurate quantification of the effect of different social restriction policies as informative references for policymakers, and it reveals that countries or regions geographically close tend to share similar pandemic evolution patterns. Second, the developed method casts the tensor completion problem in a novel embedding framework, equipped with ordinal and fusion structures, which is able to simultaneously conduct tensor completion and numerical embeddings of the countries or regions, intervention policies, and time stamps. It provides a novel treatment to tackle the informative missing issue in the COVID-19 pandemic data. Third, we establish the asymptotic consistency of the estimated data tensor, which allows all the dimensions to diverge, including the number of treated individuals, intervention policies, and time stamps. Notably, the asymptotic theory is established under an informative missing assumption, and the effect of the missing pattern has also been precisely quantified and subsequently counted into the asymptotic convergence rate. This is in sharp contrast to the missing completely at random assumption in most existing tensor completion literature (Yuan and Zhang (2016, 2017); Xia and Yuan (2019)).

The rest of the paper is organized as follows. Section 2 provides a detailed description of the COVID-19 dataset and the data processing procedure. After introducing some necessary notations in Section 3, Section 4 presents the proposed dynamic counterfactual prediction framework and its associated nonnegative tensor completion problem and develops an alternative gradient descent algorithm to tackle the resultant optimization task. Section 5

establishes the asymptotic estimation consistency of the proposed method. Application of the proposed framework on the COVID-19 pandemic data is conducted in Section 6, followed by an extensive simulation study in Section 7. Technical proofs and necessary lemmas are relegated to the Supplementary Material.

## 2. COVID-19 data.

*2.1. Data description.* We integrate data from different sources to conduct counterfactual prediction of the social mobility restriction strength on the propagation severity of the COVID-19 pandemic. Particularly, the propagation severity is measured by the daily number of death cases caused by the pandemic, as suggested in [Agarwal, Shah and Shen \(2022\)](#), which points out that the number of death cases is a more faithful measurement than the number of infections in that the testing and reporting infection results may be inconsistent across different authorities. Also, the daily number of death cases of 193 countries and regions since January 22, 2020, are made publicly available by [Dong, Du and Gardner \(2020\)](#) at <https://github.com/CSSEGISandData/COVID-19>, including the numbers on the cruise ships Diamond Princess and Ms. Zaandam. More precisely, the daily number of death cases is used to measure the propagation intensity 23 days ago, due to the delay from infection of COVID-19 to actual death. This is because, as pointed out in literature ([Lauer et al. \(2020\)](#); [Men et al. \(2023\)](#); [McAloon et al. \(2020\)](#)), the median of the incubation period of COVID-19 is about five days while the duration from symptoms onset to death is about 18 days ([Wilson et al. \(2020\)](#)).

As it is generally believed, the social mobility restriction has an essential impact on the virus' diffusion among human beings. A relatively restricted social distancing policy will generally lead to a lower infection rate and hence fewer mortality cases. The strength of the social mobility restriction across 135 countries or regions is collected by Google's mobility report at <https://www.google.com/covid19/mobility/>, which measures the daily increment or decrement of mobility trends in percentages compared to basis across six different categories of places since February 15, 2020, including retail and recreation, groceries and pharmacies, parks, transit stations, workplaces, and residential.

We extract more than 7.5 million records from Google's mobility report, where each record consists of the daily increment or decrement of mobility strength by percentiles at various categories of places in a subregion of one of the 135 countries or regions on one of the 635 days range from February 15, 2020 to November 10, 2021. Note that not every record contains the mobility changing quantities of all six categories of places. For example, one record indicates that on June 1, 2021, the mobility strength increased 7% at groceries and pharmacies, increased 2% at residential, and decreased 5% at the workplace at Abitibi Regional County Municipality, Quebec, Canada, while the mobility changing quantities at retail and recreation parks as well as transit stations are missing. In fact, a total of 5,099,553 records contain missing values for one or more categories of places, but every record consists of the mobility changing quantity for at least one place category.

*2.2. Data processing.* We first extract the cumulative numbers of death cases for all of the 193 countries or regions from January 22, 2020 to November 10, 2021, from [Dong, Du and Gardner \(2020\)](#), resulting in a  $193 \times 659$  data matrix, where its  $(i, j)$ th entry indicates the cumulative number of death cases of country or region  $i$  from January 22, 2020, to the subsequent  $j$ th day. A total of 23 countries, such as Solomon Islands, Holy See, Tonga, Palau, Kiribati, Marshall Islands, Samoa, and Micronesia, are removed from the analysis, as their cumulative numbers of death cases up to November 10, 2021, are too small for us to identify the pandemic "starting day," as in what follows. For the remaining 170 countries

or regions, we identify their first days with cumulative number of death cases more than 50 and treat those days as their corresponding starting days of the COVID-19 pandemic in the analysis. This step aligns with the intrinsic pandemic evolution of each country or region since the actual breakout time of the COVID-19 pandemic varies from country to country. Introducing a starting day could also avoid many noisy random patterns at the very beginning of the pandemic and thus enhance the signal of the data. A similar alignment of the pandemic starting day has also been employed in Agarwal, Shah and Shen (2022). We remark that the pandemic evolution may also be affected by various and possibly time-varying covariate information, such as seasonal effect, geographic location, medical treatment levels, and other social, economic, and political factors. A possible refinement of the alignment procedure shall take this covariate information into consideration, such as employing certain causal inference approaches to remove the confounding effects from the covariates.

We further delete those countries and regions with their starting day later than November 11, 2020, as we would like to study the COVID-19 pandemic evolution in a one-year period for all the studied countries, which further removes 24 countries or regions from the analysis. We then calculate the daily number of death cases for the remaining 146 countries or regions for a consecutive 365 days starting from their starting days. We finally obtain a  $146 \times 365$  daily death case matrix, whose  $(i, j)$ th entry indicates the total number of death cases of country or region  $i$  at the  $j$ th day since its starting day.

Note that the COVID-19 pandemic generally broke out locally in some major cities with relatively dense populations of a country, where international passengers frequently traveled in and out, so it is reasonable to use the same number of cumulative death cases to identify the starting date. However, shortly after the break out, the pandemic would quickly spread over the whole country, and the population density can vary drastically from one country to another. Therefore, we adjust the daily death cases by dividing the population of the corresponding country in the analysis, leading to a much more stable measure. We thus obtain a  $146 \times 365$  mortality rate matrix. The population of all the countries under investigation are available at [https://en.wikipedia.org/wiki/List\\_of\\_countries\\_and\\_dependencies\\_by\\_population](https://en.wikipedia.org/wiki/List_of_countries_and_dependencies_by_population).

Out of the 146 countries or regions, 110 of them are included in Google's mobility report, and thus we focus our analysis on these 110 countries or regions. Recall that the effects of mobility trends are reflected in the mortality rates 23 days later. Therefore, for each country or region, we extract all its records on Google's mobility report for a consecutive 365 days starting from 23 days prior to its starting day. For each record we first average the mobility changing quantities of all the available place categories in the particular day and subregion, and then for a given day in each country or region, we take the mean of the averaged mobility changing quantities of all the subregions within the country or region for that particular day. These processing steps yield a  $110 \times 365$  mobility changing quantity matrix, where its  $(i, j)$ th entry indicates the mobility changing quantity by percentile of the  $i$ th country or region in the  $j$ th day. More precisely, the positive entries indicate the increments of the human mobility trend, while the negative entries indicate the decrements.

Furthermore, a substantial chunk of the mobility data for Afghanistan and Serbia is missing in Google's mobility report, and thus we remove these two countries from the analysis. Besides, the starting day of Italy is March 2, 2020, while its mobility data are recorded only since February 15, 2020, and the starting day of Cabo Verde is September 18, 2020, while its mobility data on August 30, 2020, and September 6, 2020, are missing. To make full use of the data of these two countries, we impute the mobility flow from February 8 to February 14, 2020, in Italy, on August 30, 2020, in Cabo Verde, and on September 6, 2020, in Cabo Verde by using the data from February 15 to 28, 2020, in Italy, the data from August 25 to 29, 2020, and from August 31 to September 4, 2020, in Cabo Verde, and the data from September 1 to 5, 2020, and from September 7 to 11, 2020, in Cabo Verde, respectively. These imputations are

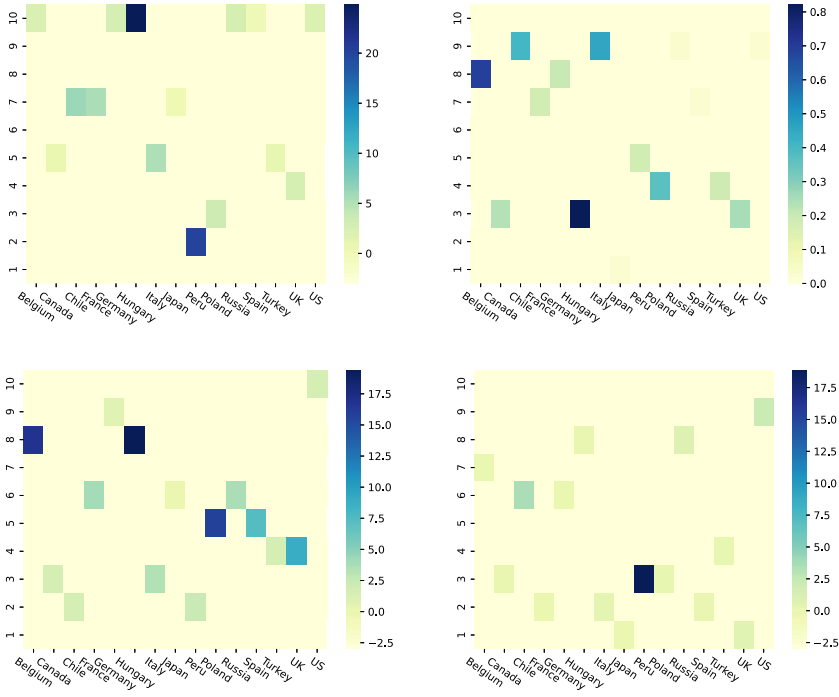


FIG. 1. Mortality rates (multiplied by  $10^6$ ) of 15 selected countries on the first, 120th, 240th, and 360th days are displayed clockwise. The horizontal axis lists the 15 countries, the vertical axis indicates the activity restriction levels, the color intensity of each block indicates each country's mortality rate based on the corresponding public policies, and the light yellow color indicates missing mortality rates.

based on the fact that the COVID-19 pandemic got more intense in the early breakout stage in Italy and Cuba Verde, and their social mobility restriction got strengthened gradually.

After all the data processing steps, we obtain two  $108 \times 365$  matrices, one about mortality rate and the other about mobility. Note that the mobility changing quantities are continuous and range from  $-46.62\%$  to  $0.93\%$ . We then evenly and sequentially cut them into 10 parts, according to the corresponding percentiles, which thus categorizes the mobility changing quantities into 10 activity restriction public policies, with policy 1 being the strongest restriction one and policy 10 the weakest restriction one since they correspond to the parts with largest and smallest human mobility decrements, respectively. Clearly, this categorization procedure ensures the same strictness of policy restriction levels across different countries. Figure 1 illustrates the mortality rates of 15 selected countries or regions in four selected days.

It is clear that each country or region will only implement one public policy to control mobility for any single day, and different countries or regions may implement different policies. Consequently, each column of the figures in Figure 1 has only one observed mortality rate, and all other entries are missing. It is thus of great interest to conduct counterfactual prediction of the mortality rates under different public policies by leveraging information from other similar countries or regions. Such predicted counterfactual outcomes can provide critical guidelines for policymakers to determine whether to elevate or delevate the public policy of social mobility control, in order to strike a better balance between COVID-19 combating and social economy development.

**3. Notations.** Throughout the paper we use boldface calligraphic Euler scripts ( $\mathcal{A}$ ) to denote tensors, boldface capital letters ( $\mathbf{A}$ ) to denote matrices, boldface lowercase letters ( $\mathbf{a}$ ) to denote vectors, and regular letters ( $a$ ) to denote scalars. For an order 3 tensor  $\mathcal{A} \in \mathbb{R}^{N \times M \times T}$ ,



$\mathcal{A}_{n,:,:} \in \mathbb{R}^{M \times T}$ ,  $\mathcal{A}_{:,m,:} \in \mathbb{R}^{N \times T}$ , and  $\mathcal{A}_{::,t} \in \mathbb{R}^{N \times M}$  are the  $n$ th horizontal,  $m$ th lateral, and  $t$ th frontal slides of  $\mathcal{A}$ , respectively, while  $\mathcal{A}_{:,m,t} \in \mathbb{R}^N$ ,  $\mathcal{A}_{n,:t} \in \mathbb{R}^M$ , and  $\mathcal{A}_{n,m,:} \in \mathbb{R}^T$  are the  $(m, t)$ th mode-1,  $(n, t)$ th mode-2, and  $(n, m)$ th mode-3 fibers of  $\mathcal{A}$ , respectively. Similarly, for a matrix  $A$ ,  $A_{i,:}$  denotes its  $i$ th row and  $A_{:,j}$  denotes its  $j$ th column. For two matrices  $A$  and  $B$ , the usual matrix product, Kronecker product, KhatriRao product (columnwise Kronecker product), and Hadamard product (entrywise product) between  $A$  and  $B$  are denoted as  $AB$ ,  $A \otimes B$ ,  $A \odot B$ , and  $A * B$ , respectively. We also use  $\mathcal{A} * \mathcal{B}$  to denote the Hadamard product between two tensors  $\mathcal{A}$  and  $\mathcal{B}$ . We use  $\|\cdot\|$ ,  $\|\cdot\|_0$ ,  $\|\cdot\|_\infty$ ,  $\|\cdot\|_{\infty,2}$ , and  $\|\cdot\|_F$  to denote the  $l_2$ -norm,  $l_0$ -norm,  $l_\infty$ -norm of a vector,  $l_{\infty,2}$ -norm of a matrix, and the Frobenius norm of a matrix or tensor, respectively. For any integer  $n$ , denote  $\mathbf{1}_n$  as the  $n$ -dimensional vector with all ones, and denote  $[n] = \{1, 2, \dots, n\}$  to be the  $n$ -set.

The mode-1 product between a tensor  $\mathcal{A} \in \mathbb{R}^{N \times M \times T}$  and a matrix  $A \in \mathbb{R}^{N' \times N}$  is denoted as  $\mathcal{A} \times_1 A \in \mathbb{R}^{N' \times M \times T}$ , whose entries are defined as

$$(\mathcal{A} \times_1 A)_{n',m,t} = \sum_{n=1}^N \mathcal{A}_{n,m,t} A_{n',n} \quad \text{for } n' \in [N], m \in [M], \text{ and } t \in [T].$$

The mode-2 or mode-3 product between a tensor and a matrix of appropriate dimension is defined in a similar fashion. Let  $\mathcal{M}_1(\mathcal{A}) \in \mathbb{R}^{N \times MT}$ ,  $\mathcal{M}_2(\mathcal{A}) \in \mathbb{R}^{M \times NT}$ , and  $\mathcal{M}_3(\mathcal{A}) \in \mathbb{R}^{T \times NM}$  be the mode-1, mode-2, and mode-3 matricization of  $\mathcal{A}$ , respectively. Precisely, the  $(m, t)$ th mode-1 fiber  $\mathcal{A}_{:,m,t}$ ,  $(n, t)$ th mode-2 fiber  $\mathcal{A}_{n,:t}$ , and  $(n, m)$ th mode-3 fiber  $\mathcal{A}_{n,m,:}$  are placed at the  $((t-1)M+m)$ th column of  $\mathcal{M}_1(\mathcal{A})$ ,  $((t-1)N+n)$ th column of  $\mathcal{M}_2(\mathcal{A})$ , and  $((m-1)N+n)$ th column of  $\mathcal{M}_3(\mathcal{A})$ , respectively, for  $n \in [N]$ ,  $m \in [M]$ , and  $t \in [T]$ .

The CANDECOMP/PARAFAC (CP) decomposition of  $\mathcal{A}$  has the form

$$(1) \quad \mathcal{A} = \sum_{k=1}^r \mathbf{u}^{(k)} \circ \mathbf{v}^{(k)} \circ \mathbf{w}^{(k)},$$

where  $\mathbf{u}^{(k)} \in \mathbb{R}^N$ ,  $\mathbf{v}^{(k)} \in \mathbb{R}^M$ , and  $\mathbf{w}^{(k)} \in \mathbb{R}^T$  for  $k \in [r]$  and  $\circ$  stands for the vector outer product. The CP-rank of a tensor  $\mathbf{u}^{(k)} \circ \mathbf{v}^{(k)} \circ \mathbf{w}^{(k)}$  is defined to be 1 (Kolda and Bader (2009)), for  $k \in [r]$ , and the CP decomposition in (1) decomposes  $\mathcal{A}$  into the summation of a number of rank-1 tensors. The minimal number of rank-1 tensors in the CP decomposition of  $\mathcal{A}$  is called the CP-rank of  $\mathcal{A}$ . Let  $\mathcal{I} \in \{0, 1\}^{r \times r \times r}$  be the order 3  $r$ -dimensional identity tensor such that  $\mathcal{I}_{i_1,i_2,i_3} = 1$  if  $i_1 = i_2 = i_3$  and 0 otherwise, and let  $U \in \mathbb{R}^{N \times r}$ ,  $V \in \mathbb{R}^{M \times r}$ , and  $W \in \mathbb{R}^{T \times r}$  such that  $U_{:,k} = \mathbf{u}^{(k)}$ ,  $V_{:,k} = \mathbf{v}^{(k)}$ , and  $W_{:,k} = \mathbf{w}^{(k)}$ . Then the CP decomposition in (1) can be equivalently rewritten as

$$\mathcal{A} = \mathcal{I} \times_1 U \times_2 V \times_3 W.$$

**4. Dynamic counterfactual prediction.** In this section we cast the COVID-19 pandemic data in a general dynamic counterfactual prediction framework in Section 4.1, then associate it with a tensor completion problem, and propose a nonnegative tensor decomposition model with various data-motivated structures in Section 4.2 to facilitate the tensor completion. An efficient alternative updated scheme is developed in Section 4.3. The key components of the proposed method in this section are illustrated in Figure 2.

**4.1. Dynamic counterfactual model.** Let  $X_{n,t} \in [M]$  be the stochastic intervention (Papadogeorgou et al. (2022)) imposed to individual  $n$  at time  $t$  and  $Y_{n,t} \in \mathbb{R}$  be the corresponding observed response affected by  $X_{n,t}$  for  $n \in [N]$  and  $t \in [T]$ , where  $N$ ,  $M$ , and  $T$  are the number of treated individuals, possible intervention policies, and time stamps under investigation, respectively. For example, in the COVID-19 pandemic data,  $X_{n,t}$  represents

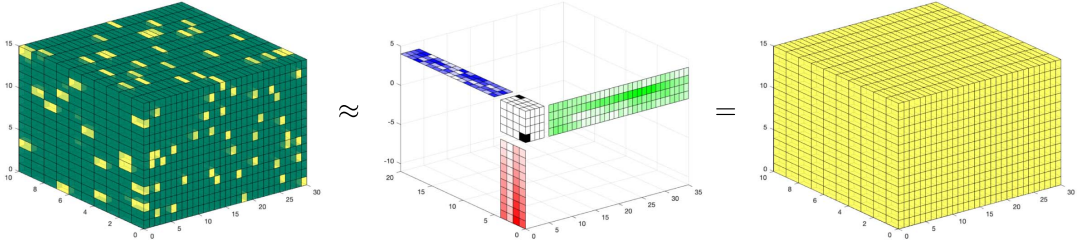


FIG. 2. The observed COVID-19 pandemic data is illustrated in the left panel, where the yellow and green entries indicate observed and potential outcomes (missing entries), respectively. The counterfactual prediction problem is cast into a nonnegative tensor CP decomposition framework with various structures, as illustrated in the middle panel. Specifically, the core tensor is an identity tensor with black and white entries representing 1 and 0, while the blue, red, and light green factor matrices are the embeddings of countries, mobility restriction policies, and time stamps, respectively. The values of these three matrices are nonnegative, where the embeddings of mobility restriction policies are ordinal (the red color gets darker), and the embeddings of time stamps possess fusion structure (some adjacent entries over time share the same light green level). Finally, after estimating the embedding matrices, we can fully estimate the outcome tensor, as illustrated by the yellow tensor in the right panel, where the estimated yellow entries can be treated as counterfactual predictions.

the human mobility restriction policy that country  $n$  adopted on the  $t$ th day, while  $Y_{n,t}$  corresponds to the exact observed mortality rate of country  $n$  on the  $t$ th day. At time  $t$  only one intervention is imposed on individual  $n$ , and thus the outcomes of other  $M - 1$  interventions on individual  $n$  can not be observed. We term these  $M - 1$  interventions as counterfactual interventions and the corresponding outcomes as counterfactuals (Höfler (2005)).

Denote  $Y_{n,t}^{(m)}$  as the potential outcome if  $X_{n,t} = m$ . We follow the standard treatment in Rubin’s potential outcome framework (Rosenbaum and Rubin (1983)) to assume the following identifiability conditions: (consistency)  $Y_{n,t} = Y_{n,t}^{(X_{n,t})}$ , (unconfoundness)  $X_{n,t}$  is independent with  $\{Y_{n,t}^{(m)}\}_{m=1}^M$ , and (strict overlap)  $P(X_{n,t} = m) \geq p_{\min}$  for some  $p_{\min} > 0$ , for  $m \in [M]$ . Clearly, the “unconfoundness” condition states that the intervention  $X_{n,t}$  is independent of the set of all potential outcomes  $\{Y_{n,t}^{(m)}\}_{m=1}^M$ , whereas the “consistent” condition states that the observed outcome  $Y_{n,t}$  strictly depends on the intervention  $X_{n,t}$ . Denote  $\theta_{n,m,t} = \mathbb{E}Y_{n,t}^{(m)}$ , the identifiability conditions allow us to write

$$(2) \quad \mathbb{E}Y_{n,t} = \sum_{m=1}^M \mathbb{E}[Y_{n,t}^{(m)} | X_{n,t} = m] P(X_{n,t} = m) = \sum_{m=1}^M \theta_{n,m,t} P(X_{n,t} = m)$$

for  $n \in [N]$ ,  $m \in [M]$  and  $t \in [T]$ .

A few remarks are in order. First, when  $X_{n,t} = m$ , the only observation is  $Y_{n,t} = Y_{n,t}^{(m)}$ , while  $Y_{n,t}^{(1)}, \dots, Y_{n,t}^{(m-1)}, Y_{n,t}^{(m+1)}, \dots, Y_{n,t}^{(M)}$  are the counterfactuals. Second, model (2) allows multiple potential interventions (Qi et al. (2020); Mo and Liu (2022); Newey and Stouli (2022)), which is a natural generalization of the conventional potential outcome framework, where only two interventions are considered (Höfler (2005); Mo, Qi and Liu (2021)), namely, control and treatment. It also allows the intervention, outcome, and its expectation to be time-varying, as in the causal framework of spatial-temporal point process (Papadogeorgou et al. (2022)). Third, the major difference distinguishes the proposed model from conventional counterfactual framework is that the expected value of the potential outcome  $\theta_{n,m,t}$  can vary from individual to individual, while existing potential outcome frameworks usually assume that individuals are randomly sampled from a source population, and thus  $Y_{1,t}^{(m)}, \dots, Y_{N,t}^{(m)}$  are independently and identically distributed for any  $m \in [M]$  and  $t \in [T]$ . This coincides with the recent set up of individualized treatment rules (Fan et al. (2022); Mo and Liu

(2022)), where individualized covariate information will provide treatment-related and possibly treatment-free effects on the potential outcomes, leading to different expected potential outcomes for different individuals, even under the same treatment at the same time.

Instead of estimating the average treatment effect of one intervention across all individuals, it is of great interest to estimate the expected responds  $\theta_{n,m,t}$ 's for different individuals. However, the heterogeneous distributions of  $Y_{n,t}^{(m)}$  cast great challenges to the estimation, since there are only  $NT$  observations, while the number of parameters to estimate is  $NMT$ , which is impossible to obtain consistent estimation unless reasonable structural assumption is imposed on  $\theta_{n,m,t}$ 's. To this end, we propose to formulate the counterfactual prediction problem in a tensor completion framework so as to reduce the intrinsic number of parameters and facilitate counterfactual prediction.

**4.2. Nonnegative tensor decomposition.** Let  $\mathcal{A} \in \mathbb{R}^{N \times M \times T}$  be the data tensor such that  $\mathcal{A}_{n,m,t} = Y_{n,t}^{(m)}$  for  $n \in [N]$ ,  $m \in [M]$ , and  $t \in [T]$ . Note that if  $X_{n,t} = m$ , then  $\mathcal{A}_{n,m,t}$  is observed while  $\mathcal{A}_{n,m',t}$  is missing for all other  $m' \neq m$ . Consequently, the problem of predicting the potential outcomes affected by the counterfactual interventions is now converted to impute the missing entries of the data tensor  $\mathcal{A}$ . To fully characterize the dynamic intervention process, we consider a nonnegative CP decomposition model,

$$(3) \quad \mathbb{E}\mathcal{A} = \Theta = \mathcal{I} \times_1 U \times_2 V \times_3 W,$$

where  $\Theta_{n,m,t} = \theta_{n,m,t}$  by the definition of  $\mathcal{A}$ ,  $\mathcal{I}$  is an order 3  $r$ -dimensional identity tensor and  $U_{n,k} \geq 0$ ,  $V_{m,k} \geq 0$ , and  $W_{t,k} \geq 0$  for  $n \in [N]$ ,  $m \in [M]$ ,  $t \in [T]$ , and  $k \in [r]$ . For more details about the tensor CP decomposition, one can refer to [Kolda and Bader \(2009\)](#) for a comprehensive review.

Clearly, (3) assures the nonnegative elements in the mean tensor  $\Theta$ , corresponding to the nonnegative expected daily mortality rates in the COVID-19 dataset. In addition, the rows of  $U$ ,  $V$ , and  $W$  serve as the numeric embeddings of the corresponding entities, interventions, and time stamps in the latent spaces. Therefore, the proposed nonnegative CP decomposition model allows one to develop efficient algorithms to conduct counterfactual prediction and numerical embedding simultaneously. Moreover, (3) reduces the inherent number of parameters in  $\Theta$  from  $NMT$  to  $(N + M + T)r$  so that consistently estimating the tensor parameter  $\Theta$  is possible based on the  $NT$  observed entries in  $\mathcal{A}$  if  $r \ll \min\{N, T\}$ . As in the tensor recommender system literature ([Bi, Qu and Shen \(2018\)](#); [Zhang et al. \(2021\)](#)) that some users may share similar preferences on items over time, the low-rank structure on  $\Theta$  essentially imposed similar evolution patterns of the causal effects with certain subgroups of treated individuals. In the COVID-19 dataset, the difference of the counterfactual pandemic evolution patterns between counties  $i$  and  $j$  is fully controlled by the difference of their embeddings  $U_{i,\cdot}$  and  $U_{j,\cdot}$ , while this heterogeneous embeddings still allow that every treated individual is sampled from a different population. Finally, we also remark that the CP decomposition (3) may not be identifiable. For instance, (3) will be still valid if one permutes the columns of  $U$ ,  $V$ , and  $W$  in the same fashion simultaneously or multiplies the  $k$ th column of  $U$ ,  $V$ ,  $W$  by positive constants  $c_U^{(k)}$ ,  $c_V^{(k)}$ , and  $c_W^{(k)}$  with  $c_U^{(k)} c_V^{(k)} c_W^{(k)} = 1$ . Fortunately, the primary goal of this paper is to estimate the expected counterfactual effect  $\Theta$ , instead of its CP decomposition factors, and hence it does not suffer from the identifiability issue induced by the tensor decomposition.

In what follows, we introduce some data-motivated properties on the CP factors to better facilitate tensor completion. Note that, in the COVID-19 dataset, the human mobility restriction policies are ordered by 10 different monotonic levels of strength. This leads to a



monotonic effect of intervention levels on the potential outcomes, following the general belief that mortality rate decreases with the strength of mobility flow restriction policies. This motivates us to introduce an ordinal structure on  $\mathbf{V}$ . Particularly, we assume that

$$\mathbf{V}_{1,:} = \Delta_1, \quad \text{and} \quad \mathbf{V}_{m,:} = \mathbf{V}_{m-1,:} + \Delta_m \quad \text{for } m = 2, \dots, M,$$

where  $\Delta_1, \Delta_2, \dots, \Delta_M$  are some  $r$ -dimensional nonnegative vectors. The expected gap between policies  $m+1$  and  $m$  for country  $n$  on the  $t$ th day is

$$\Theta_{n,m+1,t} - \Theta_{n,m,t} = \mathcal{I} \times_1 \mathbf{U}_{n,:}^\top \times_2 \Delta_{m+1}^\top \times_3 \mathbf{W}_{t,:}^\top,$$

which may vary from country to country and from day to day. Also,  $\Delta_{m+1}$  is allowed to vary with  $m \in [M-1]$  as well. It is clear that country heterogeneity and time-varying trends have been incorporated in the proposed method.

Let  $\Delta = (\Delta_1, \dots, \Delta_M)^\top$  and  $\mathbf{L} \in \{0, 1\}^{M \times M}$  be a lower triangular matrix such that  $L_{m,m'} = 1$  if  $m \geq m'$  and zero otherwise for  $m, m' \in [M]$ . It follows that  $\mathbf{V}$  admits the following decomposition:

$$(4) \quad \mathbf{V} = \mathbf{L}\Delta.$$

We remark that  $\Delta$  is the intrinsic free parameters of  $\mathbf{V}$ , and the nonnegativity constraint on  $\Delta$  is sufficient to assure both the ordinal structure and nonnegativity condition on  $\mathbf{V}$ . Therefore, the embedding  $\mathbf{V}$  is able to reflect the ordinal nature of the mobility restriction policies.

Furthermore, to incorporate the dynamic flow and smooth trend of the COVID-19 pandemic evolution, as illustrated by three selected countries in Figure 4 in Section 6, we introduce a fusion structure on  $\Theta$  through the temporal-mode factor  $\mathbf{W}$ . Specifically, let  $D: \mathbb{R}^T \rightarrow \mathbb{R}^{T-1}$  be the first-order difference operator such that  $D(\mathbf{a}) = (a_2 - a_1, \dots, a_T - a_{T-1})$  for any  $\mathbf{a} = (a_1, \dots, a_T) \in \mathbb{R}^T$ . We then introduce a fusion structure on  $\mathbf{W}$  such that

$$(5) \quad \|D(\mathbf{W}_{:,k})\|_0 \leq \kappa - 1 \quad \text{for } k \in [r],$$

where  $\kappa$  is a positive integer fusion parameter. It is clear that the constraint (5) implies that  $\mathbf{W}_{:,k}$  is successively cut into at most  $\kappa$  segments and the coordinates of  $\mathbf{W}_{:,k}$  within each segment are identical. We remark that (5) can be extended to allow different fusion parameters for different columns of  $\mathbf{W}$  at the cost of increased computational burden on model tuning. The fusion structure imposed in the temporal-mode factor matrix has been popularly employed in the dynamic tensor decomposition literature, such as dynamic tensor clustering (Sun and Li (2019)), dynamic network analysis (Zhang, Sun and Li (2020)), and dynamic tensor regression (Zhou et al. (2023)), among many others.

Let  $\Gamma_U \subset \mathbb{R}^{N \times r}$ ,  $\Gamma_V \subset \mathbb{R}^{M \times r}$ , and  $\Gamma_W(\kappa) \subset \mathbb{R}^{T \times r}$  be the domains of  $\mathbf{U}$ ,  $\mathbf{V}$ , and  $\mathbf{W}$ , respectively. Specifically,  $\Gamma_U$  is the set of all  $N \times r$  nonnegative matrices. For any  $\mathbf{V} \in \Gamma_V$ ,  $\mathbf{V}$  possesses the ordinal decomposition (4). The set  $\Gamma_W(\kappa)$  consists of all  $T \times r$  nonnegative matrices  $\mathbf{W}$  that satisfy  $\|D(\mathbf{W}_{:,k})\|_0 \leq \kappa - 1$  for  $k \in [r]$ . Denote

$$\Gamma = \{\Theta = \mathcal{I} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W} : (\mathbf{U}, \mathbf{V}, \mathbf{W}) \in \Gamma_U \times \Gamma_V \times \Gamma_W(\kappa)\}$$

as the parameter space of  $\Theta$  and  $\Omega$  as the index set of the observed entries in  $\mathcal{A}$ . The proposed nonnegative tensor completion formulation is organized as

$$(6) \quad \min_{\Theta \in \Gamma} \frac{1}{2} \|P_\Omega(\Theta - \mathcal{A})\|_F^2,$$

where  $P_\Omega(\cdot)$  is the entrywise projection that zeros the entries with indexes out of  $\Omega$ . Formally,  $P_\Omega(\Theta - \mathcal{A})$  is an  $N \times M \times T$  tensor with

$$(P_\Omega(\Theta - \mathcal{A}))_{n,m,t} = \begin{cases} \Theta_{n,m,t} - \mathcal{A}_{n,m,t} & \text{if } (n, m, t) \in \Omega; \\ 0 & \text{otherwise,} \end{cases}$$

for  $n \in [N]$ ,  $m \in [M]$  and  $t \in [T]$ . The least square formulation is computationally efficient, does not require any distributional assumption of the data tensor, as in likelihood-based formulations, and hence lets flexibility for theoretical development. The constraint set  $\Gamma$  narrows down the support of  $\Theta$  and helps avoid overfitting in terms of counterfactual prediction on the complement of  $\Omega$ .

**4.3. Alternative gradient descent algorithm.** Let  $\delta \in \{0, 1\}^{N \times M \times T}$  be the indicator tensor such that  $\delta_{n,m,t} = 1$  if  $(n, m, t) \in \Omega$  and 0 otherwise. Clearly,  $\delta_{n,m,t} = 1$  if and only if the mobility restriction policy imposed to country  $n$  on the  $t$ th day  $X_{n,t}$  equals  $m$  for  $n \in [N]$ ,  $m \in [M]$ , and  $t \in [T]$ . As a result,  $\{X_{n,m,t} : n \in [N], m \in [M], t \in [T]\}$ ,  $\Omega$ , and  $\delta$  essentially refer to the same thing but have different representation forms in different contexts. Denote the objective function as

$$f(\mathbf{U}, \mathbf{V}, \mathbf{W}) = \frac{1}{2} \|P_{\Omega}(\mathcal{I} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W} - \mathcal{A})\|_F^2 = \frac{1}{2} \|(\mathcal{I} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W} - \mathcal{A}) * \delta\|_F^2.$$

We next develop an alternative gradient descent (AGD) algorithm to minimize  $f(\mathbf{U}, \mathbf{V}, \mathbf{W})$ , subject to  $\mathbf{U} \in \Gamma_U$ ,  $\mathbf{V} \in \Gamma_V$ , and  $\mathbf{W} \in \Gamma_W$ . Clearly,  $f(\mathbf{U}, \mathbf{V}, \mathbf{W})$  is blockwise multiconvex (Xu and Yin (2013)) on  $(\mathbf{U}, \mathbf{V}, \mathbf{W})$ , yet the feasible region is nonconvex, due to the  $l_0$ -norm constraint in  $\Gamma_W$ , leading to that the optimization problem (6) can be solved only locally.

Given  $\mathbf{U}^{(s)}$ ,  $\mathbf{V}^{(s)} = \mathbf{L}\mathbf{\Delta}^{(s)}$ , and  $\mathbf{W}^{(s)}$  at iteration  $s$ , it follows from the mode-1 matricization operator that

$$\begin{aligned} f(\mathbf{U}^{(s)}, \mathbf{V}^{(s)}, \mathbf{W}^{(s)}) &= \frac{1}{2} \|(\mathbf{U}^{(s)} \mathcal{M}_1(\mathcal{I})(\mathbf{W}^{(s)} \otimes \mathbf{V}^{(s)})^\top - \mathcal{M}_1(\mathcal{A})) * \mathcal{M}_1(\delta)\|_F^2 \\ &= \frac{1}{2} \|(\mathbf{U}^{(s)}(\mathbf{W}^{(s)} \odot \mathbf{V}^{(s)})^\top) * \mathcal{M}_1(\delta) - \mathcal{M}_1(\mathcal{A} * \delta)\|_F^2. \end{aligned}$$

Therefore, the gradient of  $f(\mathbf{U}, \mathbf{V}, \mathbf{W})$ , with respect to  $\mathbf{U}$  evaluated at  $(\mathbf{U}^{(s)}, \mathbf{V}^{(s)}, \mathbf{W}^{(s)})$ , is

$$\nabla_{\mathbf{U}} f(\mathbf{U}^{(s)}, \mathbf{V}^{(s)}, \mathbf{W}^{(s)}) = \mathcal{M}_1((\mathcal{I} \times_1 \mathbf{U}^{(s)} \times_2 \mathbf{V}^{(s)} \times_3 \mathbf{W}^{(s)} - \mathcal{A}) * \delta)(\mathbf{W}^{(s)} \odot \mathbf{V}^{(s)}).$$

With a step size  $\eta > 0$ , we update  $\tilde{\mathbf{U}}^{(s+1)}$  as

$$(7) \quad \tilde{\mathbf{U}}^{(s+1)} = \mathbf{U}^{(s)} - \eta \nabla_{\mathbf{U}} f(\mathbf{U}^{(s)}, \mathbf{V}^{(s)}, \mathbf{W}^{(s)})$$

and then obtain  $\mathbf{U}^{(s+1)}$  by thresholding the negative entries of  $\tilde{\mathbf{U}}^{(s+1)}$  with zeros.

Next, given  $\mathbf{U}^{(s+1)}$  and  $\mathbf{W}^{(s)}$ , we update  $\tilde{\mathbf{\Delta}}^{(s+1)}$  along the gradient of  $f(\mathbf{U}, \mathbf{L}\mathbf{\Delta}, \mathbf{W})$ , with respect to  $\mathbf{\Delta}$  evaluated at  $(\mathbf{U}^{(s+1)}, \mathbf{L}\mathbf{\Delta}^{(s)}, \mathbf{W}^{(s)})$ , as

$$\begin{aligned} (8) \quad \tilde{\mathbf{\Delta}}^{(s+1)} &= \mathbf{\Delta}^{(s)} - \eta \nabla_{\mathbf{\Delta}} f(\mathbf{U}^{(s+1)}, \mathbf{L}\mathbf{\Delta}^{(s)}, \mathbf{W}^{(s)}) \quad \text{where} \\ \nabla_{\mathbf{\Delta}} f(\mathbf{U}^{(s+1)}, \mathbf{L}\mathbf{\Delta}^{(s)}, \mathbf{W}^{(s)}) &= \mathbf{L}^\top \mathcal{M}_2((\mathcal{I} \times_1 \mathbf{U}^{(s+1)} \times_2 \mathbf{V}^{(s)} \times_3 \mathbf{W}^{(s)} - \mathcal{A}) * \delta)(\mathbf{W}^{(s)} \odot \mathbf{U}^{(s+1)}). \end{aligned}$$

Further,  $\mathbf{\Delta}^{(s+1)}$  is obtained by thresholding the negative entries of  $\tilde{\mathbf{\Delta}}^{(s+1)}$  with zeros.

Finally, given  $\mathbf{U}^{(s+1)}$ ,  $\mathbf{V}^{(s+1)} = \mathbf{L}\mathbf{\Delta}^{(s+1)}$ , and  $\mathbf{W}^{(s)}$ , we update  $\tilde{\mathbf{W}}^{(s+1)}$  along the direction of the gradient of  $f(\mathbf{U}, \mathbf{V}, \mathbf{W})$ , with respect to  $\mathbf{W}$  evaluated at  $(\mathbf{U}^{(s+1)}, \mathbf{V}^{(s+1)}, \mathbf{W}^{(s)})$ ,

$$\begin{aligned} (9) \quad \tilde{\mathbf{W}}^{(s+1)} &= \mathbf{W}^{(s)} - \eta \nabla_{\mathbf{W}} f(\mathbf{U}^{(s+1)}, \mathbf{V}^{(s+1)}, \mathbf{W}^{(s)}) \quad \text{where} \\ \nabla_{\mathbf{W}} f(\mathbf{U}^{(s+1)}, \mathbf{V}^{(s+1)}, \mathbf{W}^{(s)}) &= \mathcal{M}_3((\mathcal{I} \times_1 \mathbf{U}^{(s+1)} \times_2 \mathbf{V}^{(s+1)} \times_3 \mathbf{W}^{(s)} - \mathcal{A}) * \delta)(\mathbf{V}^{(s+1)} \odot \mathbf{U}^{(s+1)}). \end{aligned}$$

**Algorithm 1:** Alternative gradient descent (AGD)

---

**Input** : Observed data tensor  $\mathcal{A} * \delta$ , indicator tensor  $\delta$ , CP rank  $r$ , step size  $\eta = 10^{-3}$ , and fusion parameter  $\kappa$ .

**Output:** Estimators of  $\mathbf{U}$ ,  $\mathbf{V}$  and  $\mathbf{W}$ .

- 1 Initialization. The entries of  $\mathbf{U}^{(0)}$ ,  $\mathbf{\Delta}^{(0)}$ , and  $\mathbf{W}^{(0)}$  are independently drawn from  $\text{uniform}(0, 2)$ . Set  $s = 0$ .
- 2 **repeat**
- 3   Step 1: Compute  $\tilde{\mathbf{U}}^{(s+1)}$  according to (7), and project it to the nonnegative orthant to obtain  $\mathbf{U}^{(s+1)}$ .
- 4   Step 2: Compute  $\tilde{\mathbf{\Delta}}^{(s+1)}$  according to (8), and project it to the nonnegative orthant to obtain  $\mathbf{\Delta}^{(s+1)}$ .
- 5   Step 3: Compute  $\tilde{\mathbf{W}}^{(s+1)}$  according to (9), and project it to the nonnegative orthant to obtain  $\overline{\mathbf{W}}^{(s+1)}$ . Employ the fusion operator on the columns of  $\overline{\mathbf{W}}^{(s+1)}$  to obtain  $\mathbf{W}^{(s+1)}$  as in (10).
- 6   **if**  $f(\mathbf{U}^{(s+1)}, \mathbf{V}^{(s+1)}, \mathbf{W}^{(s+1)}) \geq f(\mathbf{U}^{(s)}, \mathbf{V}^{(s)}, \mathbf{W}^{(s)})$ , **then**
- 7     | reduce step size  $\eta \leftarrow \eta/2$ , and back to line 3.
- 8   **end**
- 9    $s \leftarrow s + 1$ .
- 10 **until**  $\eta < 10^{-8}$ ;

---

Let  $\overline{\mathbf{W}}^{(s+1)}$  be the nonnegative version of  $\tilde{\mathbf{W}}^{(s+1)}$  such that the negative entries of  $\tilde{\mathbf{W}}^{(s+1)}$  are truncated by zeros. Clearly,  $\|D\overline{\mathbf{W}}^{(s+1)}\|_0 \leq \|D\tilde{\mathbf{W}}^{(s+1)}\|_0$ , implying that the fusion structure in  $\tilde{\mathbf{W}}^{(s+1)}$  will never be destroyed by the nonnegative thresholding. We then introduce a fusion operator  $F(\cdot, \cdot)$  to force the columns of  $\overline{\mathbf{W}}^{(s+1)}$  to satisfy the  $l_0$ -norm constraints (5). Particularly, for any positive integer  $d$ , positive integer  $\tau$ , and vector  $\mathbf{a} \in \mathbb{R}^d$ ,  $F(\mathbf{a}, \tau)$  first truncates  $D(\mathbf{a})$  by keeping the  $\tau - 1$  coordinates of  $D(\mathbf{a})$  with largest absolute values and zeroing all other coordinates. This partitions the coordinates of  $\mathbf{a}$  into  $\tau$  groups such that  $\mathbf{a}_{j+1}$  and  $\mathbf{a}_j$  are in the same group if the  $j$ th coordinate of the truncated version of  $D(\mathbf{a})$  is zero, for  $j \in [d - 1]$ . The  $i$ th coordinate of  $F(\mathbf{a}, \tau)$  is then defined to be the average value of all the coordinates of  $\mathbf{a}$ , inside the group that  $a_i$  belongs to, for  $i \in [d]$ . For the fusion parameter  $\kappa$ , we then obtain  $\mathbf{W}^{(s+1)}$  as

$$(10) \quad \mathbf{W}_{:,k}^{(s+1)} = F(\overline{\mathbf{W}}_{:,k}^{(s+1)}, \kappa) \quad \text{for } k \in [r].$$

The developed AGD algorithm is summarized in Algorithm 1. It employs an initial step size  $10^{-3}$  and reduces it to half only when overshooting happens. The algorithm stops when  $\eta < 10^{-8}$ , which shares the same spirit of controlling estimation error, and the estimate  $\hat{\Theta}$  is guaranteed to be close to a stationary point.

As for the CP-rank  $r$  of the approximating data tensor and the fusion parameter  $\kappa$ , we treat them as tuning parameters and select these tuning parameters through minimizing the following BIC criterion:

$$\text{BIC} = \|P_{\Omega}(\mathcal{I} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W} - \mathcal{A})\|_F^2 + (Nr + Mr + \kappa r) \log(NMT).$$

Herein, the factor  $(Nr + Mr + \kappa r)$  is the number of free parameters, which serves as the degree of freedom in the nonnegative tensor decomposition model. Similar BIC criteria have been widely employed to select tensor rank and other tuning parameters in literature

(Goutte and Amini (2010); Sun and Li (2019); Hu, Lee and Wang (2022)). We further remark that the computational complexity of the proposed method primarily lies in the computation of the gradients, while the complexities of the thresholding as well as fusion operations are negligible. Since the computation of  $\nabla_U f(U, V, W)$ ,  $\nabla_{\Delta} f(U, L\Delta, W)$ , and  $\nabla_W f(U, V, W)$  are all of order  $O(rNMT)$ , the computational complexity of Algorithm 1 is of order  $O(rNMTS)$  with  $S$  iterations, which is quite efficient and scales linearly with the size of the data tensor. Finally, we note that the scales and frequency of the negative entries in  $\tilde{U}^{(s+1)}$ ,  $\tilde{\Delta}^{(s+1)}$ , and  $\tilde{W}^{(s+1)}$  depend on the step size  $\eta$  of the algorithm. The thresholding steps for  $\tilde{U}^{(s+1)}$ ,  $\tilde{\Delta}^{(s+1)}$ , and  $\tilde{W}^{(s+1)}$  may introduce some small estimation bias. This issue can be addressed by further reducing the step size to assure positiveness of the estimated CP factors at the expense of increased computational cost.

**5. Asymptotic properties.** In this section we establish the asymptotic consistency of the estimated parameter tensor  $\hat{\Theta}$  obtained by the proposed dynamic nonnegative tensor completion method. The estimation consistency of  $\hat{\Theta}$  is evaluated via the  $\alpha$ -weighted F-norm (Mandal and Parkes (2022)) of  $\hat{\Theta} - \Theta^*$ , where  $\Theta^*$  is the underlying true parameter tensor. Specifically, let  $\alpha \in \mathbb{R}^M$  be a weighting vector with  $\alpha_m > 0$  for  $m \in [M]$ , and  $\sum_{m=1}^M \alpha_m = 1$ . For any tensor  $\mathcal{T}^{(1)} \in \mathbb{R}^{N \times M \times T}$ , its  $\alpha$ -weighted F-norm is defined as

$$\|\mathcal{T}^{(1)}\|_{F(\alpha)} = \sqrt{\langle \mathcal{T}^{(1)}, \mathcal{T}^{(1)} \rangle_{\alpha}},$$

where  $\langle \mathcal{T}^{(1)}, \mathcal{T}^{(2)} \rangle_{\alpha}$  is the  $\alpha$ -weighted inner product between two tensors  $\mathcal{T}^{(1)}, \mathcal{T}^{(2)} \in \mathbb{R}^{N \times M \times T}$ , defined as

$$\langle \mathcal{T}^{(1)}, \mathcal{T}^{(2)} \rangle_{\alpha} = \sum_{m=1}^M \alpha_m \sum_{n=1}^N \sum_{t=1}^T \mathcal{T}_{n,m,t}^{(1)} \mathcal{T}_{n,m,t}^{(2)}.$$

As  $\alpha_1, \dots, \alpha_M$  are strictly positive, the defined  $\alpha$ -weighted inner product is positive definite, symmetric, and bilinear, and hence it is indeed an inner product, inducing an valid  $\alpha$ -weighted F-norm.

The following technical assumptions are made.

**ASSUMPTION A.** Suppose that the interventions  $X_{n,t}$ ,  $n \in [N]$ ,  $t \in [T]$  are independently sampled from a multinomial distribution on  $[M]$  with parameter  $\mathbf{p} = (p_1, \dots, p_M)$ . Equivalently,  $\delta_{n,:t}$ ,  $n \in [N]$ ,  $t \in [T]$  are independently sampled from a multinomial distribution on  $\{\mathbf{e}_m\}_{m=1}^M$  with parameter  $\mathbf{p}$ , where  $\mathbf{e}_m$  is the  $m$ th canonical basis in  $\mathbb{R}^M$ .

Assumption A specifies the stochastic mechanism on the interventions with a multinomial distribution, leading to the missing mechanism in  $\mathcal{A}$ . Specifically, for each  $(n, t)$  pair, only one  $\mathcal{A}_{n,m,t}$  is observed, and all other  $\mathcal{A}_{n,m',t}$ 's with  $m' \neq m$  are missing. Clearly, the entries of  $\mathcal{A}$  are missing not at random, which differs from most existing matrix or tensor completion methods, assuming the data are missing completely at random (Yuan and Zhang (2016, 2017); Xia and Yuan (2019); Liu and Moitra (2020)) or missing at random (Xue and Qu (2021)). Moreover, the matrix or tensor regression imputation methods also fail to tackle this missing mechanism, as they often require the existence of a block of complete observations (Agarwal, Shah and Shen (2022); Bai and Ng (2021)).

**ASSUMPTION B.** Suppose that  $\mathcal{A}_{n,m,t}$ 's are independent nonnegative sub-Gaussian random variables with mean  $\Theta_{n,m,t}^*$  and a universal sub-Gaussian parameter  $\sigma > 0$ . Specifically, for any  $\lambda > 0$ ,

$$\mathbb{E} e^{\lambda(\mathcal{A}_{n,m,t} - \Theta_{n,m,t}^*)} \leq e^{\lambda^2 \sigma^2 / 2} \quad \text{for } i \in [N], m \in [M] \text{ and } t \in [T].$$

Assumption B characterizes the tail distribution behavior of the entries of  $\mathcal{A}$ , which is satisfied by many popular distributions, such as a single-side truncated normal distribution or the distribution of any nonnegative bounded random variable.

Note that a fundamental condition making tensor completion possible is the so-called incoherence condition, which essentially ensures that the rows of the factor matrices of a tensor decomposition to have comparable magnitude. Definitions of the incoherence measures vary from the tensor decomposition formats, such as the incoherence measures for orthogonal CP decomposition in Jain and Oh (2014); Cai et al. (2019), Tucker decomposition in Cai, Li and Xia (2022a), and tensor train decomposition in Cai, Li and Xia (2022b). Since, in this paper, we consider a general CP decomposition that does not require the factor matrices to be orthogonal, we introduce the following nonorthogonal CP incoherence measure for a tensor  $\Theta$  with CP-rank  $r$ :

$$(11) \quad \mu(\Theta) = \min_{(U, V, W): \Theta = \mathcal{I}_{\times 1} U \times_2 V \times_3 W} \max_{n, m, t} \{ \|U_{n,:}\|, \|V_{m,:}\|, \|W_{t,:}\| \}.$$

Clearly, though the CP decomposition of  $\Theta$  is not identifiable, the nonorthogonal CP incoherence measure is well defined, as it is defined on a best behaved decomposition in the sense that the maximum  $l_2$  norm of the rows of the factor matrices is minimized. In the following theoretical analysis, we consider those parameter  $\Theta$  has CP rank  $r$  and its nonorthogonal CP incoherence measure satisfies

$$(12) \quad \mu(\Theta) \leq \xi$$

for some  $\xi > 0$ . Clearly,  $\mu(\Theta) < \xi$  implies that  $\Theta_{n,m,t} \leq \xi^3$  for any  $n, m$ , and  $t$ . Further, for any  $W \in \Gamma_W(\kappa)$ , its  $k$ th column  $W_{:,k}$  can be split into at most  $\kappa$  segments, each consisting of identical entries. Therefore, there exists a matrix  $\tilde{w} \in [-\xi, \xi]^{\kappa \times r}$  such that the common entry value for the  $l$ th segment of  $W_{:,k}$  is  $\tilde{w}_{l,k}$ . We then call  $\tilde{w}$  a generator of  $W$ . Note that different columns may have different segment splittings.

**ASSUMPTION C.** For any  $W^{(1)}, W^{(2)} \in \Gamma_W(\kappa)$  and their corresponding generators  $\tilde{w}^{(1)}$  and  $\tilde{w}^{(2)}$ , assumes that there exists an absolute constant  $C$  such that

$$(13) \quad \|W^{(1)} - W^{(2)}\|_F \leq C\sqrt{T} \|\tilde{w}^{(1)} - \tilde{w}^{(2)}\|_{\infty, 2}.$$

Assumption C is imposed to reduce the complexity of  $\Gamma_W$ , according to the fusion structure encoded in  $W$ , such that the perturbation of  $W \in \Gamma_W$  can be well controlled by that of its generator  $\tilde{w} \in [-\xi, \xi]^{\kappa \times r}$ . It is clear that Assumption C holds with  $C = 1$  when  $\kappa = 1$  or  $\kappa = T$ . Further, as long as  $W^{(1)}$  and  $W^{(2)}$  share the same segment splittings within each column, it follows from the Hölder inequality that  $\|W^{(1)} - W^{(2)}\|_F^2 \leq T \|\tilde{w}^{(1)} - \tilde{w}^{(2)}\|_{\infty, 2}^2$ , and thus Assumption C is also satisfied with  $C = 1$ .

**THEOREM 5.1.** Under Assumptions A–C, for any estimator  $\hat{\Theta}$  satisfying

$$\frac{1}{2NT} \|(\hat{\Theta} - \mathcal{A}) * \delta\|_F^2 \leq \frac{1}{2NT} \|(\Theta^* - \mathcal{A}) * \delta\|_F^2 + \epsilon$$

with probability at least  $1 - \frac{2}{(NT)^2} - \frac{2}{(NT)^{(N+M+\kappa)r} (C\sqrt{r})^{\kappa r}}$ , it holds true that

$$\frac{1}{NT} \|\hat{\Theta} - \Theta^*\|_{F(p)}^2 \leq 4\epsilon,$$

where  $\epsilon = \max\left\{\frac{2^{11}\sigma^2((N+M+\kappa)r \log(NT) + \kappa r \log(C\sqrt{r}))}{NT}, \xi^6 \sqrt{\frac{\log(NT)}{NT}}\right\}$ .



The proof of Theorem 5.1 is provided in the Supplementary Material (Zhen and Wang (2024)). Essentially, Theorem 5.1 assures that  $\frac{1}{NT} \|\hat{\Theta} - \Theta^*\|_{F(p)}^2 = O_p(\epsilon)$  as long as the objective value of  $\hat{\Theta}$  is sufficiently close to that of  $\Theta^*$ . Herein, the convergence rate  $\epsilon$  is governed by two competing terms, where the first term depends on the number of parameters  $(N + M + \kappa)r$ , the number of observed entries  $NT$ , and the sub-Gaussian parameter  $\sigma$ , while the second term comes from the randomness of the missing pattern. We further remark that  $\xi^6$  in the second term of  $\epsilon$  is essentially an upper bound of all the squared entries in  $\Theta$ , and thus one may replace it back to have a slightly tighter upper bound. As a theoretical example with  $M, \kappa, r, \xi, \sigma$  fixed as constants, we have  $\epsilon = O(\frac{1}{T})$  if  $T = O(N)$ , and  $\epsilon = O(\frac{1}{\sqrt{NT}})$  if  $N = O(T)$ , up to some logarithm factors.

Furthermore, it follows from the definition of  $\alpha$ -weighted F-norm and Theorem 5.1 that

$$\frac{1}{NMT} \|\hat{\Theta} - \Theta^*\|_F^2 \leq \frac{4\epsilon}{Mp_{\min}}$$

with probability tending to 1, where  $p_{\min} = \min_{m \in M} p_m$ , suggesting that the convergence property of  $\hat{\Theta}$  in terms of F-norm is indeed also affected by the missing mechanism. Specifically,  $\frac{1}{NMT} \|\hat{\Theta} - \Theta^*\|_F^2$  can achieve the same convergence rate as  $\frac{1}{NT} \|\hat{\Theta} - \Theta^*\|_{F(p)}^2$  if the number of observed entries under each interventions are relatively balanced, whereas the convergence rate can be much slower if the missing pattern is unbalanced since  $p_{\min}$  can be rather small.

**6. Counterfactual prediction on COVID-19.** We now apply the proposed dynamic nonnegative tensor completion method to conduct counterfactual prediction on the COVID-19 pandemic data  $\mathcal{A} \in \mathbb{R}^{108 \times 10 \times 365}$ , where 90% of the mortality rates are missing. Before carrying out Algorithm 1, we first employ the BIC criterion to determine the CP-rank  $r$  of  $\mathbb{E}\mathcal{A}$  and the fusion parameter  $\kappa$ . Minimizing the BIC criterion over the searching space  $\{1, 2, \dots, 40\} \times \{1, 2, \dots, 40\}$  yields that  $(r, \kappa) = (18, 23)$ . Clearly, the suggested  $\kappa$  is substantially smaller than 365, suggesting that the COVID-19 pandemic indeed evolves smoothly over time.

With these tuned parameters, Algorithm 1 is carried out to produce the estimated numeric embeddings of the countries, activity restriction policies, and time stamps, denoted as  $\hat{U}$ ,  $\hat{V}$ , and  $\hat{W}$ , respectively. The counterfactual prediction on the mortality rates is then given by  $\hat{\Theta} = \mathcal{I} \times_1 \hat{U} \times_2 \hat{V} \times_3 \hat{W}$ . The submatrices corresponding to the 15 selected countries in Figure 1 of the first, 120th, 240th, and 360th mode-3 slides of  $\hat{\Theta}$  are displayed in Figure 3, which can be regarded as the estimation of the completed version of Figure 1.

It is clear from Figure 3 that the estimated mortality rate decreases with the mobility flow restriction level as the result of the ordinal structure imposed on  $V$ , which is motivated by the general belief that more strict activity restriction policy prevents large gatherings and close contacts, and is crucial for slowing down the spreads of the COVID-19 virus. More importantly, the completed tensor  $\hat{\Theta}$  also provides trustworthy prediction on the potential mortality rates for each country with different activity restriction policies. In addition, the fluctuation range of the estimation mortality rates in those countries with higher true mortality rates also tend to be larger. This is reasonable, as the pandemic evolution in those countries that have already suffered from severe pandemic situation are more sensitive to the public activity restriction policies.

To further investigate the counterfactual prediction accuracy over the temporal mode, we turn to study the mode-1 slides of  $\hat{\Theta}$ . For illustration, we focus our analysis on Brazil, Japan, and the United Kingdom. For each country the top row of Figure 4 displays its evolution curve by connecting the true mortality rates over time and a synthetic evolution curve by

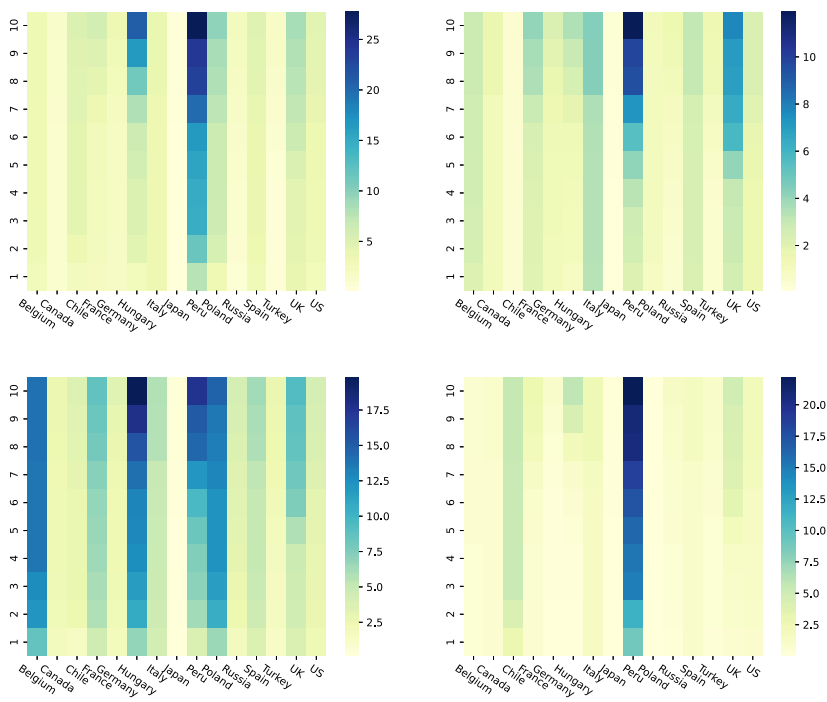


FIG. 3. The estimated mortality rates (multiplied by 10<sup>6</sup>) of the same 15 selected countries in Figure 1 on the first, 120th, 240th, and 360th days are displayed clockwise.

connecting the estimated mortality rates in  $\hat{\Theta}$ , given the actual mobility flow level every day. It is evident that the synthetic evolution trends closely coincide with the true curves in all three countries, while they are capable of avoiding overfitting by possessing smoother evolution patterns than the real curves. The bottom row of Figure 4 displays the cumulative predicted mortality rates of Brazil, Japan, and the United Kingdom over time under different

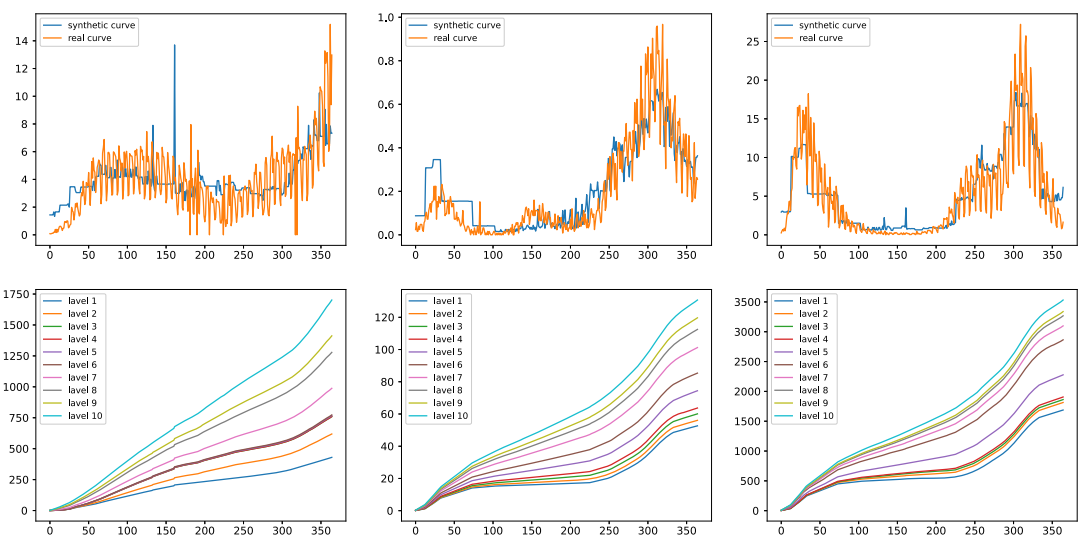


FIG. 4. The top row displays the synthetic and true evolution curves of mortality rates in Brazil (Left), Japan (Middle), and the United Kingdom (Right). The bottom row displays the predicted cumulative mortality rates by the proposed method under different social mobility restriction policies in Brazil (Left), Japan (Middle), and the United Kingdom (Right).

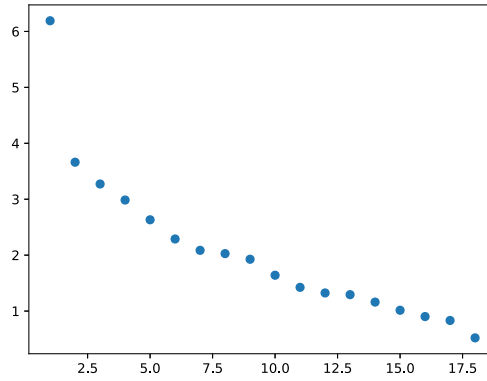


FIG. 5. A clear elbow point at the seventh leading singular value of  $\hat{U}$ .

mobility flow levels. Clearly, the predicted mortality rates under more restrictive policies are much smaller than other loose ones, and the difference among the predicted mortality rates under different policies become larger as time goes by.

As a by-product we next turn to investigate the topological structure of the estimated latent positions of the countries. The singular values of  $\hat{U}$  are displayed in Figure 5, where the seventh leading singular value is a clear elbow point (Ji and Jin (2016); Rohe, Qin and Yu (2016)), suggesting that the countries can be roughly clustered into six communities.

We hence employ a K-means algorithm with  $K = 6$  to detect communities in  $\hat{U}$ , and the results are summarized in Figure 6. It is interesting to note that community 1 consists of the United States and many Western Europe countries, such as Denmark, Estonia, Finland, France, Germany, Ireland, Italy, Portugal, Spain, United Kingdom, and Sweden. Gathering around the Caribbean Sea, Colombia, Venezuela, Panama, Nicaragua have been clustered into community 2. The majority of the countries in community 3 locate in Latin America, including Mexico, Belize, Honduras, Haiti, Ecuador, Peru, Chile, Bolivia, and Brazil. Communities 4 and 6 mainly contain countries in Eastern Europe and Africa. For example, a clique in Balkan Peninsula consists of Greece, Macedonia, Bulgaria, Romania, Hungary, Bosnia and Herzegovina, and Croatia, and a clique in southwest Sahara together with Gulf of Guinea consists of Mali, Niger, Senegal, Burkina Faso, Ghana, Cameroon, and Gabon. Community 5 only contains a few countries and is not as geographically clear as other communities, while we still can see that Argentina and Uruguay, Iraq, and Jordan are two adjacency country pairs in community 5. This community structure also suggests that the evolution of the COVID-19

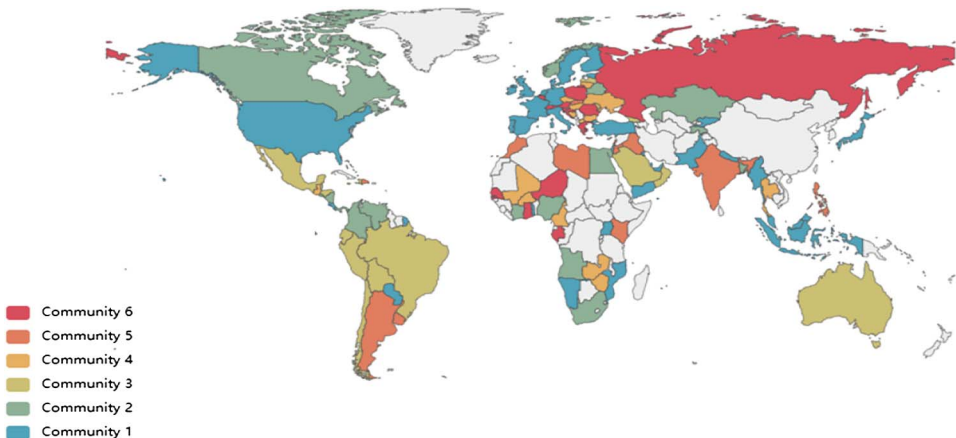


FIG. 6. Communities estimated from the COVID-19 dataset.

pandemic share great similarities among countries geographically close to each other, which is sensible, as it is easy for the virus and its variants to spread from one country to its neighboring countries.

**7. Simulation study.** To further scrutinize its empirical performance, we apply the proposed dynamic nonnegative tensor completion model, denoted as NTC, to various synthetic data and compares it against the standard tensor completion algorithm without considering the structural information (Bi, Qu and Shen (2018)), denoted as TC. We evaluate their numerical performance by the relative parameter estimation error,

$$(14) \quad \text{Err} = \frac{\|\hat{\Theta} - \Theta^*\|_F}{\|\Theta^*\|_F}.$$

**7.1. Synthetic data.** The synthetic data is generated as follows. First, we generate  $U_{n,k}^* \sim \text{Unif}(0, 2)$ ,  $V'_{m,k} \sim \text{Unif}(0, 2)$ , and  $W'_{t,k} \sim \text{Unif}(0, 2)$  for  $n \in [N]$ ,  $m \in [M]$ ,  $t \in [T]$ , and  $k \in [r]$ , where  $\text{Unif}(0, 2)$  stands for the uniform distribution in the interval  $[0, 2]$ . The  $k$ th column of  $V^*$  is then obtained by sorting the  $k$ th column of  $V'$  in ascending order for  $k \in [r]$ , and thus we have  $\Delta = L^{-1}V^*$ . For each  $k \in [r]$ , we randomly and uniformly select  $\kappa - 1$  elements from  $\{2, 3, \dots, T - 1\}$  without replacement and denote the selected values as  $\tau_1 < \dots < \tau_{\kappa-1}$ .  $W'_{t,k}$  is then computed as the average value of  $W'_{\tau_{i-1},k}, W'_{\tau_{i-1}+1,k}, \dots, W'_{\tau_i,k}$  if  $\tau_{i-1} \leq t < \tau_i$ , for  $i \in [\kappa]$ , where  $\tau_0$  and  $\tau_\kappa$  are defined as 1 and  $T + 1$ , respectively. Second, we calculate  $\Theta^* = \mathcal{I} \times_1 U^* \times_2 V^* \times_3 W^*$ , generate  $\mathcal{A} = \Theta^* + \mathcal{E}$ , and threshold the negative entries of  $\mathcal{A}$  with zeros, where the entries of  $\mathcal{E}$  are independently generated from the truncated normal distribution  $N_{[-2,2]}(0, 1)$ . Third, given  $\mathbf{p} = (p_1, \dots, p_M)^\top$ , the observation indicator tensor  $\delta$  is generated according to Assumption A:

**SCENARIO 1.** In this scenario we study the effect of the tensor size on the estimation accuracy. Specifically, we vary  $N \in \{50, 100, 150, 200, 250\}$ ,  $M \in \{10, 20, 30\}$ , and  $T \in \{100, 200, 300, 400, 500\}$  and fix  $r = 10$ ,  $\kappa = 5$ , and  $p_m = \frac{1}{M}$  for  $m \in [M]$ . For each triplet  $(N, M, T)$ , we replicate the experiment 50 times, and the averaged relative estimation errors of NTC and TC and their standard errors are reported in Tables 1–3.

**SCENARIO 2.** In this scenario we study the effect of the balance structure of the missing pattern on the estimation accuracy. Specifically, we fix  $N = 108$ ,  $M = 10$ , and  $T = 365$  to be the same of the COVID-19 dataset. We further fix  $r = 10$ ,  $\kappa = 5$ . For the weight vector  $\mathbf{p}$ , we first generate  $\mathbf{p}' \sim N_M(\mathbf{0}_M, \rho^2 \mathbf{I}_M)$  for  $\rho \in \{1, 1.5, 2, 2.5, 3\}$ , where  $\mathbf{0}_M$  is the  $M$ -dimensional vector with all zeros,  $\mathbf{I}_M$  is the  $M$ th order identity matrix, and  $N_M(\mathbf{0}_M, \rho^2 \mathbf{I}_M)$  stands for the  $M$ -variate normal distribution with mean  $\mathbf{0}_M$  and covariance matrix  $\rho^2 \mathbf{I}_M$ . We then normalize  $\mathbf{p}'$  by the soft-max operator to obtain  $\mathbf{p}$ ; that is,  $p_m = \frac{\exp(p'_m)}{\sum_{j=1}^M \exp(p'_j)}$  for  $m \in [M]$ . For each realization of  $\rho$  or equivalently  $\mathbf{p}$ , the averaged relative estimation errors of NTC and TC and their corresponding 95% confidence intervals over 50 independent replications are displayed in the left panel of Figure 7.

It is evident from Tables 1–3 and Figure 7 that NTC consistently outperforms TC in terms of averaged relative estimation errors and its standard errors in both scenarios. In addition, given  $N$  and  $M$ , the estimation error by NTC decays as  $T$  increases, whereas that of TC does not. This demonstrates the advantage of introducing a fusion structure into the time-mode factor matrix to reduce the effective degrees of free parameters, as otherwise the number of free parameters and observations will both scale linearly with  $T$ . Besides that, the estimation errors of both NTC and TC increase as  $M$  increases. This is because of the fact that, with

TABLE 1

*The averaged relative estimation errors of NTC and TC and their standard errors over 50 independent replications with  $M = 10$  in Scenario 1*

$N$		$T = 100$	$T = 200$	$T = 300$	$T = 400$	$T = 500$
50	NTC	0.0344 (0.0007)	0.0320 (0.0006)	0.0315 (0.0007)	0.0301 (0.0006)	0.0303 (0.0007)
	TC	0.0529 (0.0007)	0.0500 (0.0009)	0.0490 (0.0009)	0.0491 (0.0007)	0.0518 (0.0014)
100	NTC	0.0408 (0.0007)	0.0353 (0.0007)	0.0327 (0.0008)	0.0323 (0.0007)	0.0325 (0.0006)
	TC	0.0579 (0.0008)	0.0581 (0.0010)	0.0587 (0.0015)	0.0550 (0.0014)	0.0651 (0.0015)
150	NTC	0.0410 (0.0007)	0.0366 (0.0005)	0.0345 (0.0005)	0.0330 (0.0005)	0.0312 (0.0005)
	TC	0.0594 (0.0009)	0.0631 (0.0010)	0.0685 (0.0014)	0.0667 (0.0012)	0.0617 (0.0015)
200	NTC	0.0434 (0.0007)	0.0373 (0.0006)	0.0356 (0.0005)	0.0334 (0.0005)	0.0332 (0.0006)
	TC	0.0614 (0.0010)	0.0667 (0.0010)	0.0729 (0.0008)	0.0695 (0.0012)	0.0743 (0.0018)
250	NTC	0.0440 (0.0007)	0.0379 (0.0006)	0.0356 (0.0004)	0.0345 (0.0006)	0.0333 (0.0004)
	TC	0.0635 (0.0013)	0.0676 (0.0015)	0.0709 (0.0008)	0.0771 (0.0017)	0.0768 (0.0013)

TABLE 2

*The averaged relative estimation errors of NTC and TC and their standard errors over 50 independent replications with  $M = 20$  in Scenario 1*

$N$		$T = 100$	$T = 200$	$T = 300$	$T = 400$	$T = 500$
50	NTC	0.0388 (0.0006)	0.0334 (0.0005)	0.0318 (0.0005)	0.0298 (0.0004)	0.0289 (0.0004)
	TC	0.0590 (0.0007)	0.0558 (0.0010)	0.0623 (0.0014)	0.0610 (0.0009)	0.0761 (0.0024)
100	NTC	0.0401 (0.0008)	0.0335 (0.0005)	0.0321 (0.0005)	0.0300 (0.0003)	0.0301 (0.0004)
	TC	0.0636 (0.0009)	0.0663 (0.0007)	0.0738 (0.0016)	0.0725 (0.0005)	0.0894 (0.0024)
150	NTC	0.0420 (0.0007)	0.0354 (0.0005)	0.0319 (0.0004)	0.0305 (0.0005)	0.0286 (0.0004)
	TC	0.0700 (0.0011)	0.0742 (0.0013)	0.0791 (0.0018)	0.0853 (0.0011)	0.0822 (0.0011)
200	NTC	0.0413 (0.0007)	0.0352 (0.0006)	0.0340 (0.0005)	0.0300 (0.0005)	0.0308 (0.0006)
	TC	0.0687 (0.0010)	0.0740 (0.0009)	0.0838 (0.0013)	0.0828 (0.0004)	0.1033 (0.0019)
250	NTC	0.0443 (0.0011)	0.0407 (0.0010)	0.0326 (0.0006)	0.0375 (0.0009)	0.0317 (0.0007)
	TC	0.0796 (0.0012)	0.0853 (0.0012)	0.0807 (0.0009)	0.1009 (0.0012)	0.0908 (0.0007)

TABLE 3

*The averaged relative estimation errors of NTC and TC and their standard errors over 50 independent replications with  $M = 30$  in Scenario 1*

$N$		$T = 100$	$T = 200$	$T = 300$	$T = 400$	$T = 500$
50	NTC	0.0401 (0.0006)	0.0352 (0.0005)	0.0304 (0.0005)	0.0290 (0.0005)	0.0275 (0.0004)
	TC	0.0656 (0.0006)	0.0711 (0.0008)	0.0743 (0.0022)	0.0892 (0.0009)	0.0859 (0.0012)
100	NTC	0.0413 (0.0007)	0.0349 (0.0005)	0.0304 (0.0006)	0.0309 (0.0004)	0.0288 (0.0004)
	TC	0.0704 (0.0012)	0.0789 (0.0009)	0.0861 (0.0016)	0.0976 (0.0009)	0.0943 (0.0014)
150	NTC	0.0418 (0.0007)	0.0339 (0.0006)	0.0319 (0.0005)	0.0294 (0.0005)	0.0309 (0.0006)
	TC	0.0760 (0.0015)	0.0839 (0.0009)	0.0884 (0.0009)	0.0985 (0.0020)	0.1108 (0.0023)
200	NTC	0.0448 (0.0011)	0.0393 (0.0009)	0.0319 (0.0005)	0.0371 (0.0009)	0.0306 (0.0006)
	TC	0.0815 (0.0014)	0.0884 (0.0015)	0.0969 (0.0018)	0.1078 (0.0012)	0.1026 (0.0017)
250	NTC	0.0455 (0.0011)	0.0412 (0.0006)	0.0499 (0.0015)	0.0372 (0.0007)	0.0340 (0.0016)
	TC	0.0812 (0.0008)	0.0880 (0.0005)	0.1090 (0.0009)	0.1001 (0.0010)	0.1357 (0.0022)



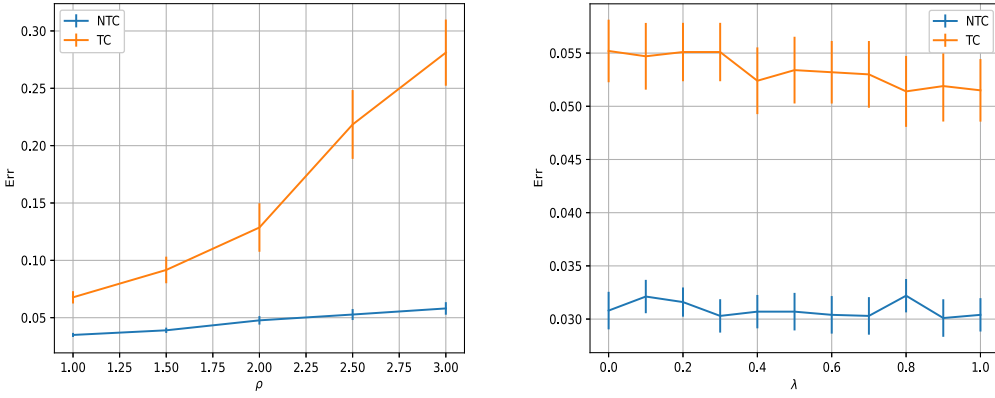


FIG. 7. *Left: The averaged relative estimation errors, given by NTC and TC, and their corresponding 95% confidence intervals for various values of  $\rho$ . Right: The averaged relative estimation errors, given by NTC and TC, and their corresponding 95% confidence intervals for various values of  $\lambda$ .*

fixing  $N$  and  $T$ , the number of parameters increases as  $M$  increases, while the number of observations remains invariant. Moreover, it can be seen from Tables 1–3 that the performance of NTC is relatively more robust than that of TC to the increasing  $M$ , and from Figure 7 that the performance of NTC is substantially more robust than that of TC to unbalance missing patterns. This is possibly due to the fact that NTC takes the ordinal structure of the intervention policies into consideration and allows for informative missing patterns, which alleviates the increment of the relative estimation error caused by the increasing  $M$  or  $\rho$ .

**7.2. Sensitivity analysis.** The proposed NTC method assumes the intervention patterns or, equivalently, the missing mechanisms to be independent among different countries and different time stamps. Yet policymakers might intend to adjust social mobility restriction policy if the mortality rates of their country or their neighboring countries significantly changed in previous days. This subsection carries out a sensitivity analysis of the proposed method to the independent intervention assumption.

Particularly, let the propensity score vector for country  $n$  on the  $t$ th day be

$$\mathbf{p}_{n,t} = (1 - \lambda)\mathbf{p}_0 + \lambda\mathbf{p}_{n,t}^{(x)},$$

where  $\mathbf{p}_0 \in \mathbb{R}^M$  is the fixed propensity score vector,  $\mathbf{p}_{n,t}^{(x)}$  is the individualized and instant propensity score vector depending on the history of the mortality rates up to the  $(t - 1)$ -day of country  $n$  and its neighboring countries, and  $\lambda \in [0, 1]$  balances the weights of  $\mathbf{p}_0$  and  $\mathbf{p}_{n,t}^{(x)}$ . We conduct the sensitivity analysis by varying  $\lambda \in \{0, 0.1, 0.2, \dots, 1.0\}$ .

We set  $N = 108$ ,  $M = 10$ , and  $T = 365$  to mimic the COVID-19 dataset. The data-generating scheme is the same as in Section 7.1, while the missing mechanism is different. Specifically, we set  $\mathbf{p}_0 = 0.1 \times \mathbf{1}_{10}$  and generate  $\mathbf{p}_{n,t}^{(x)}$  based on additional spatial-temporal information. In the first  $T_0 = 15$  days, we set  $\lambda = 0$  for all countries so that we can compute the median of the first  $NT_0$  observed outcomes, denoted as  $\gamma$ . It serves as a threshold to determine whether a mortality rate is high. For each country the closest four neighboring countries are considered. We further assume Markov property in  $\mathbf{p}_{n,t}^{(x)}$  so that  $\mathbf{p}_{n,t}^{(x)}$  depends on the mortality rate on the  $(t - 1)$ th day only. We then set  $\mathbf{p}_{n,t}^{(x)} = (\frac{1}{5}\mathbf{1}_5^\top, \mathbf{0}_5^\top)^\top$  if either the observed mortality rate of country  $n$  or the averaged observed mortality rate of its neighboring countries exceeds  $\gamma$  and set  $\mathbf{p}_{n,t}^{(x)} = (\mathbf{0}_5^\top, \frac{1}{5}\mathbf{1}_5^\top)^\top$  otherwise. For each  $\lambda$  the averaged relative estimation errors of NTC and TC and their corresponding 95% confidence intervals over 50 independent replications are displayed in the right panel of Figure 7. It is clear that both

NTC and TC are relatively robust against the variation of  $\lambda$ , suggesting that both methods do not heavily rely on the independent intervention assumption. Moreover, NTC consistently delivers smaller relative estimation error than TC, confirming the adequacy of the structural constraints.

**8. Conclusion and discussions.** This article proposes a novel nonnegative tensor completion method to conduct counterfactual prediction of mortality rates of the COVID-19 pandemic in different countries under various social mobility restriction policies. The proposed method is built upon a CP decomposition of the observed pandemic data, under the constraints that all the factor matrices are nonnegative, the factor matrix in the mobility mode is ordinal, and the factor matrix in the time mode possesses the fusion structure. The ordinal structure on the mobility mode is motivated from the general belief that restrictive social mobility policies reduce down the mortality rates. We also find out that neighboring countries that are geographically close tend to share similar pandemic evolution patterns. The effectiveness of the proposed method is also supported by numerical experiments on various synthetic datasets and asymptotic consistency in terms of parameter estimation.

We would like to remark that the independent intervention assumption is one of the key assumptions to facilitate the theoretical development, while real-life dynamic intervention mechanism can be more complicated. For example, countries may decide their intervention policy by optimizing their expected potential outcomes or taking the possibly time-varying covariate information into consideration. Consequently, it can be reasonable to allow the intervention  $X_{n,t}$  depend on the historical observations and the current available covariate information, leading to possible extension of  $\{X_{n,t}\}_{t=1}^T$  to a martingale or other discrete-time stochastic processes. Another possible extension of the current missing mechanism is to allow a small fixed part of the data tensor to be missing, such as periodic missing for one country. The current theoretic framework can still accommodate such a fixed and stochastic mixture missing pattern as long as the number of observations is substantially larger than the inherent number of parameters to be estimated. Without diluting the focus of the article, we leave there two lines of research as future works.

**Acknowledgments.** We thank the Associate Editor and two anonymous referees, whose constructive comments and suggestions have led to significant improvements of the article. We also thank Miss Ruixuan Zhao for her helpful discussion on counterfactual prediction.

**Funding.** This work is supported in part by HK RGC Grants GRF-11304520, GRF-11301521, GRF-11311022, CUHK Startup Grant 4937091, and CUHK Direct Grant 4053588.

## SUPPLEMENTARY MATERIAL

**Supplementary Materials to “Nonnegative tensor completion for dynamic counterfactual prediction on COVID-19 pandemic”** (DOI: [10.1214/23-AOAS1787SUPP](https://doi.org/10.1214/23-AOAS1787SUPP); .zip). The Supplementary Material contain necessary lemmas, technique proofs, processed data and relevant codes for all numerical experiments.

## REFERENCES

- AGARWAL, A., SHAH, D. and SHEN, D. (2022). Synthetic interventions. ArXiv preprint. Available at [arXiv:2006.07691](https://arxiv.org/abs/2006.07691).
- ATHEY, S., BAYATI, M., DOUDCHENKO, N., IMBENS, G. and KHOSRAVI, K. (2021). Matrix completion methods for causal panel data models. *J. Amer. Statist. Assoc.* **116** 1716–1730. [MR4353709 https://doi.org/10.1080/01621459.2021.1891924](https://doi.org/10.1080/01621459.2021.1891924)

- ATHEY, S., IMBENS, G. W. and WAGER, S. (2018). Approximate residual balancing: Debiased inference of average treatment effects in high dimensions. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **80** 597–623. [MR3849336 https://doi.org/10.1111/rssb.12268](https://doi.org/10.1111/rssb.12268)
- BAI, J. and NG, S. (2021). Matrix completion, counterfactuals, and factor analysis of missing data. *J. Amer. Statist. Assoc.* **116** 1746–1763. [MR4353711 https://doi.org/10.1080/01621459.2021.1967163](https://doi.org/10.1080/01621459.2021.1967163)
- BI, X., QU, A. and SHEN, X. (2018). Multilayer tensor factorization with applications to recommender systems. *Ann. Statist.* **46** 3308–3333. [MR3852653 https://doi.org/10.1214/17-AOS1659](https://doi.org/10.1214/17-AOS1659)
- CAI, C., LI, G., POOR, H. V. and CHEN, Y. (2019). Nonconvex low-rank tensor completion from noisy data. *Adv. Neural Inf. Process. Syst.* **32**.
- CAI, J.-F., LI, J. and XIA, D. (2022a). Generalized low-rank plus sparse tensor estimation by fast Riemannian optimization. *J. Amer. Statist. Assoc.* 1–17.
- CAI, J.-F., LI, J. and XIA, D. (2022b). Provable tensor-train format tensor completion by Riemannian optimization. *J. Mach. Learn. Res.* **23** Paper No. [123], 77. [MR4577075](https://doi.org/10.1080/01621459.2021.1967163)
- CHEN, B., SUN, T., ZHOU, Z. and ZENG, Y. (2019). Nonnegative tensor completion via low-rank Tucker decomposition: Model and algorithm. *IEEE Access* 95903–95914.
- DONG, E., DU, H. and GARDNER, L. (2020). An interactive web-based dashboard to track COVID-19 in real time. *Lancet Infect. Dis.* **20** 533–534. [https://doi.org/10.1016/S1473-3099\(20\)30120-1](https://doi.org/10.1016/S1473-3099(20)30120-1)
- FAN, J., MASINI, R. and MEDEIROS, M. C. (2022). Do we exploit all information for counterfactual analysis? Benefits of factor models and idiosyncratic correction. *J. Amer. Statist. Assoc.* **117** 574–590. [MR4436297 https://doi.org/10.1080/01621459.2021.2004895](https://doi.org/10.1080/01621459.2021.2004895)
- FAN, Y., LU, X., ZHAO, J., FU, H. and LIU, Y. (2022). Estimating individualized treatment rules for treatments with hierarchical structure. *Electron. J. Stat.* **16** 737–784. [MR4366820 https://doi.org/10.1214/21-ejs1948](https://doi.org/10.1214/21-ejs1948)
- FOGARTY, C. B. (2020). Studentized sensitivity analysis for the sample average treatment effect in paired observational studies. *J. Amer. Statist. Assoc.* **115** 1518–1530. [MR4143482 https://doi.org/10.1080/01621459.2019.1632072](https://doi.org/10.1080/01621459.2019.1632072)
- GOUTTE, C. and AMINI, M.-R. (2010). Probabilistic tensor factorization and model selection. In *Tensors, Kernels, and Machine Learning (TKLM 2010)* 1–4.
- HÖFLER, M. (2005). Causal inference based on counterfactuals. *BMC Med. Res. Methodol.* **5** 1–12.
- HU, J., LEE, C. and WANG, M. (2022). Generalized tensor decomposition with features on multiple modes. *J. Comput. Graph. Statist.* **31** 204–218. [MR4387221 https://doi.org/10.1080/10618600.2021.1978471](https://doi.org/10.1080/10618600.2021.1978471)
- JAIN, P. and OH, S. (2014). Provable tensor factorization with missing data. *Adv. Neural Inf. Process. Syst.* **27**.
- JI, P. and JIN, J. (2016). Coauthorship and citation networks for statisticians. *Ann. Appl. Stat.* **10** 1779–1812. [MR3592033 https://doi.org/10.1214/15-AOAS896](https://doi.org/10.1214/15-AOAS896)
- JOHANSSON, F., SHALIT, U. and SONTAQ, D. (2016). Learning representation for counterfactual inference. In *International Conference on Machine Learning* 3020–3029.
- KOLDA, T. G. and BADER, B. W. (2009). Tensor decompositions and applications. *SIAM Rev.* **51** 455–500. [MR2535056 https://doi.org/10.1137/07070111X](https://doi.org/10.1137/07070111X)
- LAUER, S. A., GRANTZ, K. H., BI, Q., JONES, F. K., ZHENG, Q., MEREDITH, H. R., AZMAN, A. S., REICH, N. G. and LESSLER, J. (2020). The incubation period of coronavirus disease 2019 (COVID-19) from publicly reported confirmed cases: Estimation and application. *Ann. Intern. Med.* **172** 577–582. <https://doi.org/10.7326/M20-0504>
- LIU, A. and MOITRA, A. (2020). Tensor completion made practical. In *Advances in Neural Information Processing Systems* **33** 18905–8916.
- MANDAL, D. and PARKES, D. (2022). Weighted tensor completion for time-series causal inference. ArXiv preprint. Available at [arXiv:1902.04646](https://arxiv.org/abs/1902.04646).
- MICALOON, C., COLLINS, Á., HUNT, K., BARBER, A., BYRNE, A. W., CASEY, M., GRIFFIN, J., LANE, E., MCEVOY, D. et al. (2020). Incubation period of COVID-19: A rapid systematic review and meta-analysis of observational research. *BMJ Open* **10** e039652.
- MEN, K., LI, Y., WANG, X., ZHANG, G., HU, J., GAO, Y., HAN, A. and LIU, W. (2023). Estimate the incubation period of coronavirus 2019 (COVID-19). *Comput. Biol. Med.* **158** 106794.
- MO, W. and LIU, Y. (2022). Efficient learning of optimal individualized treatment rules for heteroscedastic or misspecified treatment-free effect models. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **84** 440–472. [MR4412993 https://doi.org/10.1111/rssb.12474](https://doi.org/10.1111/rssb.12474)
- MO, W., QI, Z. and LIU, Y. (2021). Learning optimal distributionally robust individualized treatment rules. *J. Amer. Statist. Assoc.* **116** 659–674. [MR4270012 https://doi.org/10.1080/01621459.2020.1796359](https://doi.org/10.1080/01621459.2020.1796359)
- NEWAY, W. K. and STOULI, S. (2022). Heterogeneous coefficients, control variables and identification of multiple treatment effects. *Biometrika* **109** 865–872. [MR4472853 https://doi.org/10.1093/biomet/asab060](https://doi.org/10.1093/biomet/asab060)
- PAPADOGEORGOU, G., IMAI, K., LYALL, J. and LI, F. (2022). Causal inference with spatio-temporal data: Estimating the effects of airstrikes on insurgent violence in Iraq. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **84** 1969–1999. [MR4515563 https://doi.org/10.1111/rssb.12548](https://doi.org/10.1111/rssb.12548)

- PEARL, J. (2009). Causal inference in statistics: An overview. *Stat. Surv.* **3** 96–146. [MR2545291](https://doi.org/10.1214/09-SS057) <https://doi.org/10.1214/09-SS057>
- POULOS, J., ALBANESE, A., MERCATANTI, A. and LI, F. (2022). Retrospective causal inference via matrix completion, with an evaluation of the effect of European integration on cross-border employment. ArXiv preprint. Available at [arXiv:2106.00788](https://arxiv.org/abs/2106.00788).
- QI, Z., LIU, D., FU, H. and LIU, Y. (2020). Multi-armed angle-based direct learning for estimating optimal individualized treatment rules with various outcomes. *J. Amer. Statist. Assoc.* **115** 678–691. [MR4107672](https://doi.org/10.1080/01621459.2018.1529597) <https://doi.org/10.1080/01621459.2018.1529597>
- QUICK, C., DEY, R. and LIN, X. (2021). Regression models for understanding COVID-19 epidemic dynamics with incomplete data. *J. Amer. Statist. Assoc.* **116** 1561–1577. [MR4353694](https://doi.org/10.1080/01621459.2021.2001339) <https://doi.org/10.1080/01621459.2021.2001339>
- ROHE, K., QIN, T. and YU, B. (2016). Co-clustering directed graphs to discover asymmetries and directional communities. *Proc. Natl. Acad. Sci. USA* **113** 12679–12684. [MR3576189](https://doi.org/10.1073/pnas.1525793113) <https://doi.org/10.1073/pnas.1525793113>
- ROSENBAUM, P. R. and RUBIN, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika* **70** 41–55. [MR0742974](https://doi.org/10.1093/biomet/70.1.41) <https://doi.org/10.1093/biomet/70.1.41>
- SINHA, S. and CHAKRABORTY, M. (2022). Causal analysis and prediction of human mobility in the US during the COVID-19 pandemic. ArXiv preprint. Available at [arXiv:2111.12272](https://arxiv.org/abs/2111.12272).
- SUN, W. W. and LI, L. (2019). Dynamic tensor clustering. *J. Amer. Statist. Assoc.* **114** 1894–1907. [MR4047308](https://doi.org/10.1080/01621459.2018.1527701) <https://doi.org/10.1080/01621459.2018.1527701>
- WILSON, N., KVALSVIG, A., BARNARD, L. T. and BAKER, M. G. (2020). Case-fatality risk estimates for COVID-19 calculated by using a lag time for fatality. *Emerg. Infect. Dis.* **26** 1339–1441. <https://doi.org/10.3201/eid2606.200320>
- XIA, D. and YUAN, M. (2019). On polynomial time methods for exact low-rank tensor completion. *Found. Comput. Math.* **19** 1265–1313. [MR4029842](https://doi.org/10.1007/s10208-018-09408-6) <https://doi.org/10.1007/s10208-018-09408-6>
- XU, Y. and YIN, W. (2013). A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM J. Imaging Sci.* **6** 1758–1789. [MR3105787](https://doi.org/10.1137/120887795) <https://doi.org/10.1137/120887795>
- XUE, F. and QU, A. (2021). Integrating multisource block-wise missing data in model selection. *J. Amer. Statist. Assoc.* **116** 1914–1927. [MR4353722](https://doi.org/10.1080/01621459.2020.1751176) <https://doi.org/10.1080/01621459.2020.1751176>
- YADLOWSKY, S., PELLEGRINI, F., LIONETTO, F., BRAUNE, S. and TIAN, L. (2021). Estimation and validation of ratio-based conditional average treatment effects using observational data. *J. Amer. Statist. Assoc.* **116** 335–352. [MR4227698](https://doi.org/10.1080/01621459.2020.1772080) <https://doi.org/10.1080/01621459.2020.1772080>
- YAN, H., ZHU, Y., GU, J., HUANG, Y., SUN, H., ZHANG, X., WANG, Y., QIU, Y. and CHEN, S. X. (2021). Better strategies for containing COVID-19 pandemic: A study of 25 countries via a vSIADR model. *Proc. R. Soc. A* **477** Paper No. 20200440, 25. [MR4258333](https://doi.org/10.1098/rspa.2020.0440)
- YE, Y., ZHANG, Q., CAO, Z., CHEN, F. Y., YAN, H., STANLEY, H. E. and ZENG, D. D. (2021). Impact of export restrictions on the global personal protective equipment trade network during COVID-19. *Adv. Theory Simul.* 2100352.
- YUAN, M. and ZHANG, C.-H. (2016). On tensor completion via nuclear norm minimization. *Found. Comput. Math.* **16** 1031–1068. [MR3529132](https://doi.org/10.1007/s10208-015-9269-5) <https://doi.org/10.1007/s10208-015-9269-5>
- YUAN, M. and ZHANG, C.-H. (2017). Incoherent tensor norms and their applications in higher order tensor completion. *IEEE Trans. Inf. Theory* **63** 6753–6766. [MR3707566](https://doi.org/10.1109/TIT.2017.2724549) <https://doi.org/10.1109/TIT.2017.2724549>
- ZHANG, J., SUN, W. W. and LI, L. (2020). Mixed-effect time-varying network model and application in brain connectivity analysis. *J. Amer. Statist. Assoc.* **115** 2022–2036. [MR4189774](https://doi.org/10.1080/01621459.2019.1677242) <https://doi.org/10.1080/01621459.2019.1677242>
- ZHANG, X. and NG, M. K. (2022). Sparse nonnegative tensor factorization and completion with noisy observations. *IEEE Trans. Inf. Theory* **68** 2551–2572. [MR4413569](https://doi.org/10.1109/TIT.2021.3104356)
- ZHANG, Y., BI, X., TANG, N. and QU, A. (2021). Dynamic tensor recommender systems. *J. Mach. Learn. Res.* **22** Paper No. 65, 35. [MR4253758](https://doi.org/10.1080/01621459.2021.1938082)
- ZHEN, Y. and WANG, J. (2024). Supplement to “Nonnegative tensor completion for dynamic counterfactual prediction on COVID-19 pandemic.” <https://doi.org/10.1214/23-AOAS1787SUPP>
- ZHOU, J., SUN, W. W., ZHANG, J. and LI, L. (2023). Partially observed dynamic tensor response regression. *J. Amer. Statist. Assoc.* **118** 424–439. [MR4571132](https://doi.org/10.1080/01621459.2021.1938082) <https://doi.org/10.1080/01621459.2021.1938082>