

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

A Bayesian multivariate factor analysis model for evaluating an intervention by using observational time series data on multiple outcomes

This is the author's manuscript

Original Citation:

Availability:

This version is available <http://hdl.handle.net/2318/1754986> since 2020-11-19T14:40:32Z

Published version:

DOI:10.1111/rssa.12569

Terms of use:

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)

A Bayesian multivariate factor analysis model for evaluating an intervention using observational time-series data on multiple outcomes

Pantelis Samartsidis¹, Shaun R. Seaman¹, Silvia Montagna², André Charlett³, Matthew Hickman⁴ and Daniela De Angelis¹

¹*MRC Biostatistics Unit, University of Cambridge, University Forvie Site, Robinson Way, Cambridge CB2 0SR, UK* ²*Dipartimento di Scienze Economico-sociali e Matematico-statistiche (ESOMAS), University of Torino, Corso Unione Sovietica, 218 bis, Torino 10134, Italy* ³*Statistics, Modelling and Economics Department, Public Health England, 61 Colindale Avenue, London NW9 5EQ, UK* ⁴*School of Social and Community Medicine, University of Bristol, Oakfield Grove, Clifton, Bristol BS8 2BN, UK*

March 6, 2020

Abstract

A problem frequently encountered in many areas of scientific research is that of estimating the impact of a non-randomised binary intervention on an outcome of interest using time-series data on units that received the intervention ('treated') and units that did not ('controls'). One popular estimation method in this setting is based on the factor analysis (FA) model. The FA model is fitted to the pre-intervention outcome data on treated units and all the outcome data on control units, and the counterfactual treatment-free post-intervention outcomes of the former are predicted from the fitted model. Intervention effects are estimated as the observed outcomes minus these predicted counterfactual outcomes. In this article, we propose a model that extends the FA model for estimating intervention effects by: 1) jointly modelling the multiple outcomes to exploit shared variability, and 2) assuming an autoregressive structure on factors to account for temporal correlations in the outcome. Using simulation studies, we show that the proposed method can improve the precision of the intervention effect estimates and achieve better control of Type I error rate (compared to the FA model), especially when either the number of pre-intervention measurements or the number of control units is small. We apply our method to estimate the impact of stricter alcohol licensing policies on alcohol-related harms.

1 Introduction

In this work, we consider the problem of estimating the causal effect of an intervention on an outcome of interest in the setting where: i) the intervention is binary; ii) assignment of the sample units to the intervention is non-randomised; iii) only a small number of units are treated; and iv) there are multiple measurements of the outcome both before and after the intervention occurs.

This problem is frequently encountered in various fields of scientific research, including econometrics, epidemiology, marketing, public health and political science. For example: [Card \(1990\)](#) studied the effect that the mass migration in 1980 of Cubans to Miami had on the Miami's labour market, by treating Miami as having received the 'intervention' of mass Cuban migration and comparing it with other US states not subject to such migrations; [Cavallo *et al.* \(2013\)](#) assessed the impact that large-scale natural disasters, such as earthquakes and storms, had on the gross domestic product of a country by comparing countries subject to such natural disasters (the 'intervention') with countries not experiencing such disasters; [de Vocht \(2016\)](#) investigated whether the increased use of mobile phones (the 'intervention') led to an increase in incidences of certain types of brain cancer in England.

A general difficulty when estimating the causal effect of an intervention from observational data is the potential existence of confounding variables. These are variables which affect both the outcome of

interest and the probability of being assigned to the intervention. Failure to account for confounding can lead to biased estimation of causal effects. When the number of units receiving the intervention is large, propensity score methods (Robins *et al.*, 2000) can be used. However, when few units receive the intervention, there is not enough information to fit propensity score models. For this reason, several new methodologies for causal inference in the setting where i)-iv) apply have recently been proposed. For a recent review, see Samartsidis *et al.* (2019).

Many of these methods, including Abadie *et al.* (2010), Hsiao *et al.* (2012), Gobillon and Magnac (2016), Chan and Kwok (2016) and Xu (2017), build upon the *factor analysis* (FA) model. FA is a natural way to adjust for confounding. The model allows for unobserved confounders that remain constant over time but have a time-varying effect on the outcome. However, current methodologies based on the FA model have shortcomings. Firstly, they can only be applied to a single outcome at a time. When there is more than one correlated outcome, it might be more efficient to model them jointly. Secondly, none of the aforementioned methods explicitly models the temporal correlation between the multiple measurements of the outcome. Modelling autocorrelation may improve efficiency. Thirdly, when the total number of units is small, it is hard to perform inference for the causal effects using the existing methods. Finally, some of the approaches above require specifying the number of factors in the model. Although guidance is provided for how to use the data to choose this number, inference using these methods does not account for this data-dependent choice and hence tends to be anti-conservative.

In this paper, we attempt to address these shortcomings. We consider extensions of the FA model that can exploit the correlation between different outcomes and the temporal correlation within each outcome, leading to more efficient estimates. Also, by taking a Bayesian approach, we can obtain credible intervals for the causal effects that account for the uncertainty in the number of factors. We contribute to the literature on causal inference in the setting where i)-iv) apply, in three ways. Firstly, we develop a novel approach that uses multivariate outcomes. An alternative multivariate model is suggested by Robbins *et al.* (2017). However, their method is designed for high-dimensional data and its utility in a small-data setting is unclear. Secondly, our method is one of the few that model temporal correlation within each outcome. Brodersen *et al.* (2015) recently proposed the *Causal Impact* method to account for such correlation. However, *Causal Impact* can only be applied to a single treated unit and single outcome at a time. Thirdly, our use of the Bayesian approach allows a more inference on the causal effects of interest in comparison to other FA-based approaches.

Our method has connections with various applications of FA in contexts other than causal inference. More specifically, this is not the first time that a multivariate factor model has been used in practice. De Vito *et al.* (2016, 2018) and Avalos-Pacheco *et al.* (2018) demonstrate the benefits of taking a multivariate approach in genomic applications when dealing with multiple studies rather than multiple outcomes. Assuming a temporal structure in the FA model is also not uncommon; see, for example, McAlinn *et al.* (2017) for a recent application in macroeconomics. However, to our knowledge, this is the first time that either of these extensions (joint outcome modelling and explicit modelling of the temporal correlation) to the standard FA model has been implemented in a causal inference problem and the benefits of using them demonstrated. Our methodology, together with other causal inference methodologies based on FA, is related to the *latent class analysis* causal inference approach (Lanza *et al.*, 2013; Bartolucci *et al.*, 2016; Tullio and Bartolucci, 2019, LCA). In LCA, there is a fixed number of classes (the analogue of factors in FA) and the distribution of the outcomes on each unit and at each time point depends on the unobserved class of that unit at that time point. Hence, both LCA and FA attempt to model the variability in the outcome using latent variables. A difference between LCA and FA is that classes are discrete whereas factors are continuous. Despite being similar in spirit to FA approaches, LCA-based causal inference methods cannot be used in the setting where i)-iv) apply, mainly because they require estimation of the propensity score, which is problematic when the number of treated units is small. Moreover, these methods focus on the causal effect of the intervention on the probability that an individual belongs to a certain class, whereas in our problem the interest is in the effect of the intervention directly on the outcomes.

The paper is structured as follows. Section 2 introduces our motivating example. Section 3.1 introduces the notation and causal framework. The standard FA model is formulated in Section 3.2. Section 3.3 presents the proposed methodology. Section 3.4 describes how our method accounts for the uncertainty regarding the true number of factors. Prior distributions and posterior sampling are discussed in Section 3.5. Section 3.6 describes how point estimates and inferences are obtained for the causal effects of interest. In Section 4 we perform a series of simulation studies to evaluate the utility of the proposed methodology, using the standard FA model as our benchmark. Section 5 describes the application of methods to our motivating dataset. Finally, Section 6 contains a discussion and suggests some possible

directions for future research.

The data that are analysed in the paper and the programs that were used to analyse them can be obtained from www.url.co.uk.

2 Motivating example

Alcohol consumption has an adverse impact on society, being responsible for a number of harmful health conditions and behaviours. National policy makers have long focused on the development of effective strategies to limit these negative effects. For example, the 2003 Licensing Act¹ in England and Wales enables local authorities to develop *cumulative impact policies* (CIPs) that is, automatically reject new licensing applications unless these are supported by evidence that grant will not negatively impact on surrounding premises.

In a recent study, [de Vocht et al. \(2017\)](#) assessed the impact that CIPs had on alcohol-related harms. They collected data on four alcohol-related outcomes: hospital admission rate per 10000 people; violent crimes rate per 1000 people; sexual crime rate per 1000 people; and antisocial behaviour incidence rate per 1000 people. The data on each outcome were collected quarterly for the period mid-2009 to 2015. The authors defined intervention sites as local councils implementing a CIP in 2012, and control sites as local councils that did not adopt a CIP at any time during the study period. They identified 5 treated and 86 control sites in England and Wales.

In Section 5, we demonstrate our proposed methodology using a subset of data in [de Vocht et al. \(2017\)](#). We exclude data from one treated site (Tyneside) because the intervention was implemented earlier in this site and from 9 control sites because of missing values in some of the outcomes. Finally, we only use data up to mid-2013, because trends in the following months might be due to changes in the way crimes were reported ([de Vocht et al., 2017](#)). Figure 1 shows the data.

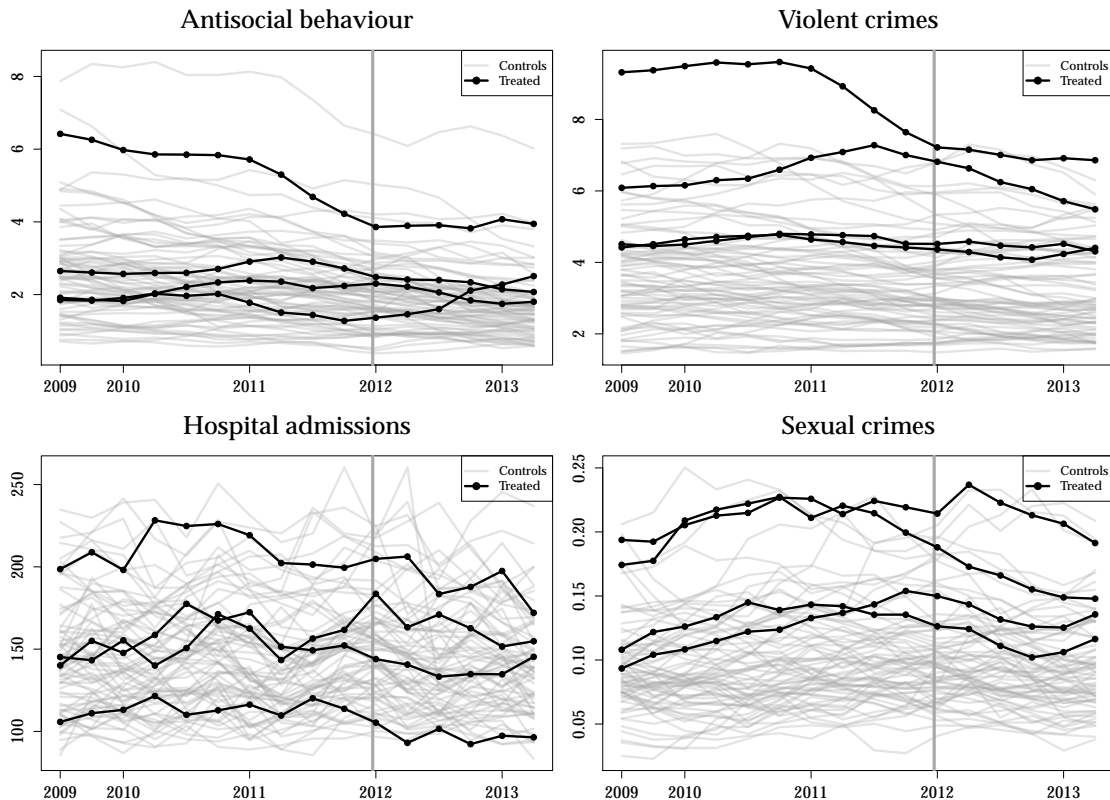


Figure 1: The alcohol licensing dataset. Each plot shows the time-series for each unit for one of the four outcomes. Control and treated councils are represented in grey and black colour, respectively. The vertical grey line indicates the introduction of the intervention.

¹<http://www.legislation.gov.uk/ukpga/2003/17/contents>

3 Model specifications

3.1 Notation and causal framework

We have observations y_{itk} , where $i = 1, \dots, n$ indexes the units, $t = 1, \dots, T$ indexes the time points and $k = 1, \dots, K$ indexes the outcomes. The units are ordered so that the first n_1 are the controls, i.e. units that do not receive the intervention during the course of the study. For the remaining $n_2 = n - n_1$ units, there is a time point T_1 after which they all receive the intervention. We refer to these units as the *treated* units. Let r_i be a binary indicator of whether unit i is treated. The study can be split into two periods, the pre-intervention period consisting of the first T_1 time points when none of the n units has the intervention, and the post-intervention period consisting of the remaining $T_2 = T - T_1$ time points when the intervention is in place for the n_1 treated units.

In this paper, we adopt the *Rubin causal model* (Rubin, 1974; Holland, 1986). This means that for each treated unit i ($i > n_1$), time t after intervention (i.e. $t > T_1$) and outcome k there are two potential outcomes $y_{itk}^{(0)}$ and $y_{itk}^{(1)}$; $y_{itk}^{(0)}$ represents the outcome that would have been observed if the intervention had not been applied and $y_{itk}^{(1)}$ is the outcome that would be observed if the intervention were applied. Hence, the causal effect of the intervention for any $i > n_1$, $t > T_1$ and k is given by

$$\theta_{itk} = y_{itk}^{(1)} - y_{itk}^{(0)}. \quad (1)$$

We are further interested in the average treatment effect on outcome k at time t in the treated units, ϑ_{tk} , defined as

$$\vartheta_{tk} = \frac{1}{n_2} \sum_{i=n_1+1}^n \theta_{itk}. \quad (2)$$

For treated units before intervention (i.e. $i > n_1$ and $t \leq T_1$) and for control units at all times (i.e. $i \leq n_1$ and all t), $y_{itk}^{(0)} = y_{itk}$ for all k , and so is observed. We do not observe $y_{itk}^{(0)}$ for treated units after intervention, and so causal effects (1) and (2) are not observed. Our approach is to assume a model for $y_{itk}^{(0)}$, use this to obtain predictions $\hat{y}_{itk}^{(0)}$ for the counterfactuals $y_{itk}^{(0)}$ for $i > n_1$ and $t > T_1$, and then estimate (1) and (2) as $\hat{\theta}_{itk} = y_{itk}^{(1)} - \hat{y}_{itk}^{(0)}$ and $\hat{\vartheta}_{tk} = \frac{1}{n_2} \sum_{i=n_1+1}^n \hat{\theta}_{itk}$, respectively.

3.2 The factor analysis model for a single outcome

For time-series observational data, a model frequently used for $y_{itk}^{(0)}$ is the FA model, also known as the *interactive fixed effects* model (Bai, 2009). Gobillon and Magnac (2016), Chan and Kwok (2016) and Xu (2017) use the FA model for causal inference in the setting we are investigating, i.e. that where i)-iv) apply. Abadie *et al.* (2010) and Hsiao *et al.* (2012) show that their proposed estimators of the counterfactuals $y_{itk}^{(0)}$ ($i > n_1$ and $t > T_1$) are unbiased when the FA model is the data-generating mechanism.

The FA model for the k -th outcome assumes that

$$y_{itk}^{(0)} = \boldsymbol{\gamma}_{ik}^\top \mathbf{s}_{tk} + \varepsilon_{itk}, \quad (3)$$

where $\boldsymbol{\gamma}_{ik} = (\gamma_{ik1}, \dots, \gamma_{ikp_1})^\top$ is the p_1 -vector of unobserved unit-specific loadings for outcome k , $\mathbf{s}_{tk} = (s_{tk1}, \dots, s_{tkp_1})^\top \sim N_{p_1}(\mathbf{0}, \mathbf{I})$ is the p_1 -vector of unobserved time-specific factors for outcome k , and $\varepsilon_{itk} \sim N(0, \psi_{ik}^2)$ is the error term. One can view the loadings $\boldsymbol{\gamma}_{ik}$ as unit characteristics that remain constant over time and the factors \mathbf{s}_{tk} as their time-varying effect on the potential outcome. Xu (2017) refers to \mathbf{s}_{tk} as ‘shocks’; $\boldsymbol{\gamma}_{ik}$ describe the magnitude of the effect that these shocks have on unit i ’s outcome k . Variables that are predictive of $y_{itk}^{(0)}$ but are not affected by the intervention can be incorporated as covariates \mathbf{x}_{it} in the FA model by replacing Eqn. (3) with

$$y_{itk}^{(0)} = \boldsymbol{\gamma}_{ik}^\top \mathbf{s}_{tk} + \boldsymbol{\beta}_k^\top \mathbf{x}_{it} + \varepsilon_{itk}. \quad (4)$$

Such variables may include observed confounders measured before time T_1 . For simplicity, we shall omit such covariates until Section 3.5.

We note in passing that a special case of the FA model is the *difference-in-differences* (DID) model (Angrist and Pischke, 2009; Jones and Rice, 2011). This is the FA model with $p_1 = 2$, $s_{tk1} = 1$ and $\gamma_{ik2} = 1$, i.e. fixed effects for units and time points. The DID model assumes that the ‘shocks’ at each time point affect all the units in the same way. This model is frequently used for causal inference with time-series observational data.

Recall that the potential outcome $y_{itk}^{(0)}$ is observed at all times (i.e. $t = 1, \dots, T$) for control units ($r_i = 0$) but at only the pre-intervention times (i.e. $t = 1, \dots, T_1$) for treated units ($r_i = 1$). Model (3) is fitted to these observed data, considering the post-intervention outcomes $y_{itk}^{(0)}$ on treated units as missing. The resulting estimator $\hat{\theta}_{itk}$ of the intervention effect on unit i and outcome k at time t is asymptotically unbiased as $n_1 \rightarrow \infty$ and $T_1 \rightarrow \infty$, provided that r_i is independent of $\varepsilon_{i1k}, \dots, \varepsilon_{iT_k}$ (and assuming regularity conditions) (Xu, 2017).

There are two intuitive ways to understand why this asymptotic unbiasedness holds. First, as n_1 and T_1 become larger (assuming fixed $n - n_1$ and $T - T_1$), the amount of data for learning about the factors \mathbf{s}_{tk} and the loadings $\boldsymbol{\gamma}_{ik}$ increases, so that the factors and loadings (and hence the expectation of $y_{itk}^{(0)}$) are increasingly accurately estimated. Second, by letting $\mathbf{y}_{tk}^{(0)} = \left(y_{1tk}^{(0)}, \dots, y_{ntk}^{(0)} \right)^\top$, defining $\boldsymbol{\Gamma}_k$ as the $n \times p_1$ matrix with i -th row $\boldsymbol{\gamma}_{ik}$, and letting $\boldsymbol{\epsilon}_{tk} = (\epsilon_{1tk}, \dots, \epsilon_{ntk})^\top$, Equation (3) implies that

$$\mathbf{y}_{tk}^{(0)} = \boldsymbol{\Gamma}_k \mathbf{s}_{tk} + \boldsymbol{\epsilon}_{tk}, \quad (5)$$

for all t and k . From this, marginally, i.e. integrating out the factors and error terms, we have that

$$\text{Cov}\left(\mathbf{y}_{tk}^{(0)}\right) = \boldsymbol{\Gamma}_k \boldsymbol{\Gamma}_k^\top + \boldsymbol{\Psi}_k, \quad (6)$$

where $\boldsymbol{\Psi}_k = \text{diag}\{\psi_{1k}^2, \dots, \psi_{nk}^2\}$. Hence, the FA model assumes that the covariance of the potential (treatment-free) outcomes of the n units is the same at all time points. The pre-intervention data are used to learn about this covariance which is then used to predict the (counterfactual) potential outcomes of the treated units after intervention from the (observed) potential outcomes of the control units after intervention. The larger are n_1 and T_1 , the more information is available to estimate $\boldsymbol{\Gamma}_k$ and $\boldsymbol{\Psi}_k$, and hence the more accurately we can estimate them (and, from them, the expectation of $y_{itk}^{(0)}$).

It is worth noting that the FA model allows for a certain form of unmeasured confounding. This is because the aforementioned asymptotic unbiasedness property of $\hat{\theta}_{itk}$ does not require r_i to be independent of $\boldsymbol{\gamma}_{ik}$. If $\boldsymbol{\gamma}_{ik}$ is indeed associated with r_i then, because it is also associated with $y_{itk}^{(0)}$ (see Eqn. (3)), it is an (unobserved) confounder.

3.3 Extending the FA model

Our proposed model involves two extensions to the FA model: joint outcome modelling and temporal dependence. We present these two extensions separately, although the model we finally propose, the ‘MVFA+AR’ model, includes both extensions.

Joint outcome modelling. The classical FA model considers each of the K different outcomes independently; it makes no assumptions about correlations between outcomes k and k' ($k' \neq k$). In situations where the different outcomes are measures of, or are influenced by, a common underlying process, this may be an inefficient way to estimate intervention effects. For example, the outcomes GDP and employment rate can be considered to be two measures of the underlying health of an economy; and rates of hospital admission, violent crime, sexual crime and antisocial behaviour are all influenced by problematic alcohol use. In these situations, part of the variability of the different outcomes is shared. Such shared variability can be modelled using a multivariate FA (MVFA) model. As we explain below, the MVFA model enables the counterfactual post-intervention k -th outcomes of the treated units to be estimated using the data on all K outcomes, rather than (as in the FA model) just the data on the k -th outcome. This makes it possible to estimate these counterfactual outcomes — and hence the intervention effects — more precisely.

The MVFA model assumes that

$$y_{itk}^{(0)} = \boldsymbol{\gamma}_{ik}^\top \mathbf{s}_{tk} + \boldsymbol{\lambda}_i^\top \mathbf{f}_{tk} + \varepsilon_{itk}, \quad (7)$$

where $\boldsymbol{\gamma}_{ik}$ and \mathbf{s}_{tk} are as defined earlier (i.e. they are unit-specific loadings and time-specific factors, both of which are specific to the k -th outcome), $\boldsymbol{\lambda}_i$ is the p_2 -vector of unit-specific loadings that are shared across outcomes, $\mathbf{f}_{tk} \sim N_{p_2}(\mathbf{0}, \mathbf{I})$ is a p_2 -vector of time-specific factors for $\boldsymbol{\lambda}_i$, and $\varepsilon_{itk} \sim N(0, \psi_{ik}^2)$ is the error term. Again, covariates \mathbf{x}_{it} can be included in the model by adding the term $\boldsymbol{\beta}_k^\top \mathbf{x}_{it}$ to the right hand side of Equation (7).

The interpretation of the MVFA model follows that of the FA model. More specifically, as well as $\boldsymbol{\gamma}_{ik}$, we now have $\boldsymbol{\lambda}_i$, which can be thought of as unit-specific unobserved variables that affect all outcomes;

their effect on outcome k at time t is quantified by the factor \mathbf{f}_{tk} . One way to think about the benefit of jointly modelling the outcomes when estimating the counterfactual outcomes is by appreciating that the joint model learns about $\boldsymbol{\lambda}_i$, the unit-specific loadings common to all the outcomes, from the data on all the outcomes. This means that $\boldsymbol{\gamma}_{ik}^\top \mathbf{s}_{tk} + \boldsymbol{\lambda}_i^\top \mathbf{f}_{tk}$, the expectation of the counterfactual k -th outcome, is more accurately estimated than in the case when modelling the outcomes independently. A second, alternative way to think about this benefit is to consider the covariance matrix for $y_{ik}^{(0)}$. For the FA model, this is given by Equation (6). For the MVFA model, we have that

$$\text{Cov}\left(\mathbf{y}_{ik}^{(0)}\right) = \boldsymbol{\Gamma}_k \boldsymbol{\Gamma}_k^\top + \boldsymbol{\Lambda} \boldsymbol{\Lambda}^\top + \boldsymbol{\Psi}_k,$$

for each k and t , where $\boldsymbol{\Lambda}$ is the $n \times p_2$ matrix with i -th row $\boldsymbol{\lambda}_i$. By modelling the outcomes jointly, this covariance can be estimated more accurately, because the part attributable to the shared factors, i.e. $\boldsymbol{\Lambda} \boldsymbol{\Lambda}^\top$, is estimated using data from all the K outcomes. Since this covariance is used to predict the (counterfactual) potential outcomes of the treated units after the intervention, estimating it more accurately should lead to more accurate estimation of those outcomes.

We expect the benefit of jointly modelling the outcomes to be greatest when T_1 is small. In this situation, there is little data to learn the unit-specific loadings, and so the gain from learning about some of them (specifically $\boldsymbol{\lambda}_i$) using all the outcomes is likely to be most marked. Also, the greater is the proportion of factors that are common, i.e. the larger the p_2/p_1 , the greater is likely to be the benefit from using the MVFA model. Note that joint modelling of multiple outcomes should be beneficial in terms of improving the precision of the estimate of the causal effect even when the effect of intervention on only one of the K outcomes is of interest. Also note that time-dependent variables that are predictive of $y_{itk}^{(0)}$ but which are affected by the intervention cannot be included as covariates in the MVFA model (as in Equation (4)). However, they can be used as additional outcomes in the MVFA model.

Modelling temporal dependence. The effect of the unit-specific loading $\boldsymbol{\gamma}_{ik}$ (or $\boldsymbol{\lambda}_i$) on the outcomes at times t and t' is represented by \mathbf{s}_{tk} and $\mathbf{s}_{t'k}$ (or \mathbf{f}_{tk} and $\mathbf{f}_{t'k}$). It may be reasonable to believe that this effect is likely to be more similar at two nearby times than at two distant times, e.g. \mathbf{s}_{tk} and $\mathbf{s}_{t+1,k}$ are likely to be more similar than \mathbf{s}_{tk} and $\mathbf{s}_{t+10,k}$. Neither the FA nor MVFA model described above takes this time ordering into account. We can take into account the time ordering by assuming that for each outcome k , the factors are generated by an AR(1) process. Specifically, we assume that, for each k and $j = 1, \dots, p_1 + p_2$ we have that

$$s_{tkj} = \rho_{kj} s_{t-1,kj} + \eta_{tkj}, \quad (8)$$

where $s_{tkj} = f_{tk,j-p_1}$ for $p_1 < j \leq p_1 + p_2$, $\rho_{kj} \in (-1, 1)$ are persistent parameters, and $\eta_{tkj} \sim \text{N}(0, 1)$.

Assuming that factors are generated by an AR(1) process may improve prediction of the counterfactual outcomes $y_{itk}^{(0)}$ ($i > n_1, t > T_1$) and hence increase the precision of the intervention effect estimates. This can become clear as follows. By integrating out the factors and error terms, we find that, for $t' \neq t$,

$$\text{Cov}\left(\mathbf{y}_{tk}^{(0)}, \mathbf{y}_{t'k}^{(0)}\right) = \boldsymbol{\Gamma}_k \text{Cov}\left(\mathbf{s}_{tk}, \mathbf{s}_{t'k}\right) \boldsymbol{\Gamma}_k^\top + \boldsymbol{\Lambda} \text{Cov}\left(\mathbf{f}_{tk}, \mathbf{f}_{t'k}\right) \boldsymbol{\Lambda}^\top. \quad (9)$$

Equation (9) shows that by assuming an AR(1) prior for the factors, an *a priori* correlation both between y_{itk} and $y_{it'k}$ is allowed, as well as between y_{itk} and $y_{i't'k}$, where $i \neq i'$. If these correlations are strong, the sharing of information across time points can lead to more accurate estimates of the counterfactuals. This does not happen in the standard FA model which assumes that $\rho_{kj} = 0$, and therefore the right hand side of Eq. 9 reduces to $\mathbf{0}$.

We expect that assuming an AR structure for the factors will increase efficiency in settings where n_1 is small and T_1 large. In these settings, there are few observations per time point and therefore factors cannot be estimated accurately. By assuming an AR(1) structure, we allow for the sharing of information between nearby times points. When T_1 is small, there may be less advantage, because there is then less information to estimate ρ_{kj} .

We call the FA model with this AR(1) structure the ‘FA+AR’ model and the MVFA model with AR(1) structure the ‘MVFA+AR’ model.

3.4 Choosing the number of factors

One of the challenges when implementing FA is choosing the total number of factors in the model. Many authors have proposed solutions for this problem; for example, Bai and Ng (2002) propose some criteria to choose the number of factors; Lopes and West (2004) develop a reversible jump Markov chain

Monte Carlo (MCMC) algorithm to estimate the number of factors; [Carvalho et al. \(2008\)](#) take an evolutionary stochastic model search approach; [Srivastava et al. \(2017\)](#) use a continuous shrinkage prior on the loadings. In this work, we account for the uncertainty in p_1 and p_2 by assuming a *multiplicative Gamma shrinkage process* (MGPS) prior ([Bhattacharya and Dunson, 2011](#)) on the loadings.

We use the MGPS for both the outcome-specific loadings $\boldsymbol{\gamma}_{ik}$ and shared loadings $\boldsymbol{\lambda}_i$. We shall describe how this prior works for the outcome-specific loadings; for the shared loadings, the specifications are analogous. Let the loading vector $\boldsymbol{\gamma}_{ik}$ be of dimension $p_1 = \infty$. The MGPS assumes that for each $j = 1, \dots, \infty$ we have

$$\gamma_{ikj} \sim \text{N}\left(0, \frac{1}{\phi_{ikj}\tau_{kj}}\right), \quad (10)$$

where ϕ_{ikj} and τ_{kj} (both > 0) are the local and global shrinkage parameters, respectively, such that

$$\tau_{kj} = \prod_{\ell=1}^j \delta_{k\ell}. \quad (11)$$

For appropriately chosen priors on ϕ_{ikj} and δ_{kj} , the product $\phi_{ikj}\tau_{kj}$ increases, thus encouraging the magnitude of the elements of $\boldsymbol{\gamma}_{ik}$ to decrease progressively towards zero. Hence, although the number of columns in each matrix $\boldsymbol{\Gamma}_k$ is infinite there will be a column such that all columns after this column have an L_1 norm of almost zero, indicating that no more factors are required for the dataset under consideration.

In practice, it is not possible to carry out computations when loadings are infinite-dimensional. So, we let $\boldsymbol{\gamma}_{ik}$ be of dimension k_1 (and k_2 for the shared loadings), where k_1 is sufficiently large. This approach can be computationally wasteful when p_1 is much smaller than the specified k_1 . However, one can easily detect this through a pilot run of the algorithm used to simulate from the posterior; if most of the columns of $\boldsymbol{\Gamma}_{ik}$ have an L_1 norm that is very low, then it is recommended to decrease k_1 in the final run. Alternatively, an adaptive way to determine k_1 is discussed by [Bhattacharya and Dunson \(2011\)](#).

The MGPS prior can be used to perform inference on the number of factors. At iteration ℓ of MCMC (which we use to draw samples from the posterior), let $d^{(\ell)}$ denote the total number of columns in $\boldsymbol{\Gamma}_k$ whose absolute elements $|\gamma_{1kj}|, \dots, |\gamma_{nkj}|$ are all below a pre-specified threshold m . The effective number of factors at iteration ℓ is $k_1 - d^{(\ell)}$. Therefore, one can use the posterior distribution of $k_1 - d^{(\ell)}$ to estimate the total number of factors in the model (e.g. as the posterior median of this distribution) and construct credible intervals (using the quantiles). This approach is sensitive to the choice of threshold m .

The reasons that we use the MGPS prior instead of the other methods that we mention are twofold. Firstly, the MGPS allows for a conjugate formulation of the model, which simplifies posterior sampling. Secondly, it has been shown that this method performs well in a wide range of applications (e.g. [Montagna et al. \(2012, 2018a,b\)](#)).

3.5 Prior distributions and MCMC algorithm

Prior distributions are as follows. For all i and k , we let the variance parameters $\psi_{ik}^2 \sim \text{InverseGamma}(0.001, 0.001)$. For the AR parameters, $\rho_{kj} \sim \text{Uniform}(-1, 1)$ for all k and j . For the shrinkage parameters, we follow recommendations by [Bhattacharya and Dunson \(2011\)](#) and let $\phi_{ikj} \sim \text{Gamma}(\frac{3}{2}, \frac{3}{2})$ for all i, k and j , $\delta_{k1} \sim \text{Gamma}(2.1, 1)$ for all k , and $\delta_{kj} \sim \text{Gamma}(3.1, 1)$ for $j > 1$. If covariates \boldsymbol{x}_{it} are included in the MVFA+AR model, we let the regression coefficients $\boldsymbol{\beta}_k \sim \text{N}(0, 10^3 \boldsymbol{I})$ for all k .

The posterior distribution resulting from the MVFA+AR model of Equations (7), (8), (10) and (11) and the prior distributions stated this section is analytically intractable. We therefore use MCMC to draw samples from it. In particular, we propose a hybrid Gibbs sampler where each parameter (or block of parameters) is sampled from its full conditional given the remaining parameters, using either Gibbs or Metropolis-Hastings steps. The main challenge is to simulate from high-dimensional normal full conditionals. More specifically, the vector of factors $\boldsymbol{f}_k = \left(\boldsymbol{s}_{1k}^\top, \boldsymbol{f}_{1k}^\top, \dots, \boldsymbol{s}_{Tk}^\top, \boldsymbol{f}_{Tk}^\top\right)^\top$ for each outcome is drawn from a $T(k_1 + k_2)$ -dimensional normal distribution, and the vector of loadings $\tilde{\boldsymbol{\lambda}}_i = \left(\boldsymbol{\lambda}_i^\top, \boldsymbol{\gamma}_{i1}^\top, \dots, \boldsymbol{\gamma}_{iK}^\top\right)^\top$ for each unit is drawn from a $(Kk_1 + k_2)$ -dimensional normal distribution. We perform both these updates with good computational efficiency using the method of [Rue \(2001\)](#). The update of AR hyperparameters ρ_{jk} is also challenging, because these parameters have bounded support and it is not possible to simulate them directly from their full conditionals. To overcome this issue, we update ρ_{jk} with Metropolis-Hastings steps using the proposals developed by [Kastner and Frühwirth-Schnatter \(2014\)](#). The remaining model parameters $\boldsymbol{\beta}_k$, ϕ_{ikj} , δ_{kj} and ψ_{ik}^2 can be easily drawn from their

full conditionals. For full details of the MCMC algorithm, see section A of the web-based supplementary material, where we also provide a sketch of the sampler. Similar MCMC algorithms can be used to draw from the posterior distribution of the FA, FA+AR and MVFA models.

We emphasise that the factors and loadings are not identifiable. Since we are not interested in interpreting these parameters but only in the counterfactuals (which are identifiable), we choose not to impose any identifiability constraints. Users interested in interpreting these parameters can resort to one of the existing approaches for ensuring identifiability, see e.g. Section 12.1.3 of [Murphy \(2012\)](#) for a fairly recent overview. One method is to restrict the loading matrix to the class of lower diagonal matrices ([Geweke and Zhou, 1996](#)).

3.6 Point estimation and inference

Samples from the posterior distribution are used to obtain samples from the posterior distribution of the causal effect θ_{itk} . First, we simulate from the posterior predictive distribution of the counterfactual outcomes

$$y_{itk}^{(0,\ell)} = (\boldsymbol{\gamma}_{ik}^{(\ell)})^\top \mathbf{s}_{tk}^{(\ell)} + (\boldsymbol{\lambda}_i^{(\ell)})^\top \mathbf{f}_{tk}^{(\ell)} + \varepsilon_{itk}^{(\ell)}, \quad (i > n_1, t > T_1) \quad (12)$$

where $\varepsilon_{itk}^{(\ell)} \sim N(0, (\psi_{ik}^2)^{(\ell)})$ and ℓ indexes the MCMC draw. Then, samples $\theta_{itk}^{(\ell)}$ from the posterior of the individual effect θ_{itk} are readily available as $\theta_{itk}^{(\ell)} = y_{itk}^{(0,\ell)} - y_{itk}^{(0)}$. Similarly, samples $\vartheta_{tk}^{(\ell)}$ from the posterior of the average treatment effect ϑ_{tk} are obtained as $\vartheta_{tk}^{(\ell)} = \frac{1}{n_2} \sum_{i=n_1+1}^n \theta_{itk}^{(\ell)}$. We can use these to calculate point estimates and perform inference. For instance, the point estimate of ϑ_{tk} will be $\frac{1}{L} \sum_{\ell=1}^L \vartheta_{tk}^{(\ell)}$, where L is the number of MCMC samples and the 95% credible interval for ϑ_{tk} will be given by the 2.5% and 97.5% percentiles of the $\vartheta_{tk}^{(\ell)}$. To test for a positive intervention effect, one can estimate the posterior probability that $\vartheta_{tk} > 0$ as $\frac{1}{L} \sum_{\ell=1}^L \mathbb{I}(\vartheta_{tk}^{(\ell)} > 0)$, where $\mathbb{I}(\cdot)$ is the indicator function.

4 Simulation studies

4.1 Setting

We performed a series of simulation studies to answer the question of whether one can obtain estimates of ϑ_{tk} that are more precise than those obtained from the standard FA model by: (i) modelling multiple outcomes jointly; (ii) assuming an AR(1) structure for the factors; and (iii) doing both simultaneously.

Each dataset (from a total of 10000) was simulated as follows. We used the MVFA+AR model to generate data on $n = 35$ units with $T_1 = 40$ pre-intervention and $T_2 = 5$ post-intervention time points. There were $K = 3$ outcomes, $p_1 = 2$ outcome-specific loadings, $p_2 = 4$ shared loadings, and the persistent parameters of the factors were $\rho_{kj} = 0.9$ for all k and j ². s_{0kj} were drawn from a $N(0, 1)$ distribution. For each k, i and j , we drew the loadings from a $N(0, 1)$. Finally, for each k and i , we set $\psi_{ik}^2 = 1/3$.

We randomly chose $n_2 = 5$ treated units from these 35 units using the expected values on the first outcome. To introduce unobserved confounding, each unit had selection probability proportional to

$$\text{expit} \left\{ \kappa \sum_{t=41}^{45} \left[\boldsymbol{\gamma}_{i1}^\top \mathbf{s}_{t1} + \boldsymbol{\lambda}_i^\top \mathbf{f}_{t1} \right] \right\},$$

of being selected to be a treated unit, where $\text{expit}(\cdot) = \exp(\cdot)/(1 + \exp(\cdot))$. The value of κ controls the degree of unobserved confounding; $\kappa = 0$ means no unobserved confounding, and $\kappa > 0$ means that units with larger expected values of the post-intervention (possibly counterfactual) treatment-free first outcome are more likely to be treated. We chose $\kappa = 0.75$ because we found that a simple t-test comparing $\mathcal{Y}_1 = \{y_{1,41,1}^{(0)}, \dots, y_{30,41,1}^{(0)}\}$ and $\mathcal{Y}_2 = \{y_{31,41,1}^{(0)}, \dots, y_{35,41,1}^{(0)}\}$ had roughly 17.5% rejection rate (this would be around 5% if the elements of \mathcal{Y}_1 and \mathcal{Y}_2 were exchangeable). This procedure gave us datasets of $n_1 = 30$ control units, $n_2 = 5$ treated units (Setup I), and $T_1 = 40$ and $T_2 = 5$.

We expected the answers to questions (i)-(iii) to depend on T_1 and n_1 . To obtain datasets with fewer than 40 pre-intervention time points and/or fewer than 30 control units, we discarded the data in the first $40 - T_1$ pre-intervention time points and/or randomly discarded $30 - n_1$ control units. The values of (T_1, n_1) in setups II-IX are (40, 30), (40, 15) and (40, 5), (20, 30), (20, 15), (20, 5), (10, 30), (10, 15)

²We set the variance of the error terms η_{tkj} in Eq. (8) to $1/(1 - \rho_{jk})$ for all k and j , so that $s_{tkj} \sim N(0, 1)$ for all t, j and k .

and (10, 5), respectively. The total number of treated units ($n_2 = 5$) and post-intervention observations ($T_2 = 5$) were common to all setups.

The point estimate of the ϑ_{tk} is

$$\hat{\vartheta}_{tk} = \frac{1}{n_2} \sum_{i=n_1+1}^n \left[y_{itk}^{(1)} - \hat{y}_{itk}^{(0)} \right] = \vartheta_{tk} + \frac{1}{n_2} \sum_{i=n_1+1}^n \left[y_{itk}^{(0)} - \hat{y}_{itk}^{(0)} \right].$$

So, if the average (over simulated datasets and over treated units) value of $\hat{y}_{itk}^{(0)}$ ($i > n_1$ and $t > T_1$) is equal to the average (over simulated datasets and over treated units) of $y_{itk}^{(0)}$, then $\hat{\vartheta}_{tk}$ will be unbiased for any value of θ_{itk} . Similarly, if the credible interval for $\sum_{i=n_1+1}^n y_{itk}^{(0)}$ contains the true value of this sum, then the credible interval for ϑ_{tk} will also contain the true value of the ϑ_{tk} for any θ_{itk} . Thus, it sufficed to study the case where $\theta_{itk} = 0$.

We fit the following models to all datasets: FA, FA+AR, MVFA and MVFA+AR. Note that all these models were correctly specified for the data that we generated. For all methods, we ran MCMC for 31250 iterations, applied a thinning factor of 25 to obtain a total of 1250 posterior draws, discarded the first 250 as a burn-in and used the remaining 1000 for inference. Models FA and FA+AR are designed for univariate outcomes and hence we applied these to each of the outcomes in turn. For all models, we used the MGPS prior, setting $k_1 = k_2 = 12$. We used the posterior mean as the point estimate of the ϑ_{tk} ; credible intervals were obtained using the 2.5% and 97.5% quantiles of the posterior distribution.

We compared the performance of the models in terms of the bias and standard error of the point estimates for the ϑ_{tk} 's, and mean width and false positive rate of the 95% credible intervals of the ϑ_{tk} 's. As a measure of 'power', we used the probability of detecting an intervention effect when the true ϑ_{tk} was not equal to zero. Here, we defined 'detection' as a credible interval that excludes zero. For convenience, we assumed that θ_{itk} was the same for all units, times and outcomes.

4.2 Results

The results for the first outcome ($k = 1$) are summarised in Table 1 and Figures 2 and 3. Table 1 presents bias, standard error, mean credible interval width and false positive rate for $(k, t) = (1, T_1 + 1)$ and $(k, t) = (1, T)$. Figures 2 and 3 show power for $(k, t) = (1, T_1 + 1)$ and $(k, t) = (1, T)$, respectively. In Section B of the web-based supplementary material, we present results for $(k, t) = (2, T_1 + 1)$ (Table 1 and Figure 1), and for $(k, t) = (3, T)$ (Table 1 and Figure 2).

To answer question (i), we compare the results obtained from the MVFA+AR and MVFA models to the results obtained from the FA+AR and FA models, respectively. In settings where $T_1 < 40$ and $n_1 \geq 15$ (i.e. Setups IV, V, VII and VIII), we see that joint modelling of outcomes leads to considerable gains in precision: the MVFA+AR and MVFA models decrease the standard error of the point estimates and the mean credible interval width in these settings (see Table 1 and Table 1 of the web-based supplementary material Section B). The gains in efficiency are also apparent from the power: Figures 2-3 and 1-2 of web-based supplementary material Section B show that the use of a multivariate model instead of a univariate model can substantially improve the chance of detecting an intervention effect. For example, for $(k, t) = (1, T)$ and Setup VII, we find that when $\vartheta_{T_1+1,1} = 1.2$, the intervention effect is detected by the MVFA model with probability 78% whereas it is detected by the FA model with probability 66%. The power curves for $(k, T) = (2, T_1 + 1)$ and $(k, t) = (3, T)$ (web-based supplementary material Section B) reveal a similar pattern. We find no settings in which the univariate models outperform the corresponding multivariate models for any of the performance measures that we consider.

To answer question (ii), we compare the results obtained from the MVFA+AR and FA+AR models to those obtained from the MVFA and FA models, respectively. We find that the inclusion of the AR component leads to either improved or unchanged performance. The improvements occur mainly in settings where $n_1 = 5$ (i.e. Setups III, VI and IX). In these settings with few control units, the FA and MVFA models perform very poorly in terms of bias and false positive rate for outcome $k = 1$ (see Table 1). The FA+AR and MVFA+AR models offer significant improvements in terms of both bias and false positive rate compared to the FA and MVFA models. For example, for $(k, t) = (1, T)$ and Setup VI, the false positive rate is 8.9% for the FA model whereas it is 5.7% for the MVFA+AR model. Note that for outcomes $k = 2$ and $k = 3$, the bias and false positive rate in Setups III and VI are not as high as for outcome $k = 1$ (see Table 1 of web-based supplementary material Section B). The reason is that we have chosen the treated units using the expected outcomes on $k = 1$ and therefore the effect of confounding is greater for $k = 1$. The inclusion of the AR component also leads to gains in efficiency in Setups III and VI, as it reduces both the standard error of the point estimates and the mean credible

interval width (Table 1). The gains in power can be moderate. For example, for $(k, t) = (1, T)$ and Setup III, the intervention effect is detected with probability 90% by the FA+AR model and 83% by the FA model when $\vartheta_{T_1+1,1} = 1.5$ (Figure 2). The improvement in power is more prominent for outcome $k = 2$ (because the bias of all methods is close to zero for this outcome). For instance, in Setup III, a $\vartheta_{T_1+1,3} = 1.5$ is detected with probability 86% by the FA+AR model and 72% by the FA model (Figure 1 of web-based supplementary material Section B).

We find no setups in which the MVFA+AR model performs better than both the MVFA and FA+AR models. The reason is that, as we explained earlier in Section 3.3, the two proposed extensions (joint outcome modelling and the AR(1) prior on factors) improve upon FA in very different settings: the former when T_1 is small and n_1 is large; the latter when T_1 is large and n_1 is small. On the other hand, we find no settings where either FA+AR or MVFA model outperforms the MVFA+AR model. Therefore, the advantage of the MVFA+AR model is that it can be used in all settings. For this reason, we suggest that this is the model that should be used in practice.

The efficiency gains obtained by using either the FA+AR, MVFA and MVFA+AR models will depend not only on T_1 and n_1 but also on the total number of outcomes K , the number of factors $p_1 + p_2$ and the ratio p_1/p_2 . As K increases, λ_i will be estimated with higher precision. As the total number of factors $p_1 + p_2$ increases, the total number of model parameters that need to be estimated increases. Hence, sharing of information (either between outcomes using the MVFA/MVFA+AR models or between time points using the FA+AR/MVFA+AR models) is more important for larger values of $p_1 + p_2$. Finally, the MVFA and MVFA+AR models are well suited to applications where the ratio p_1/p_2 is low i.e. where the number of shared loadings is large compared to the number of outcome-specific loadings.

5 Application to the motivating dataset

5.1 Model details

In this section we apply the proposed methodology to the alcohol licensing dataset introduced in Section 2. The dataset consists of data on $K = 4$ outcomes relating to the harms of alcohol consumption in society. For each outcome, there are $n_1 = 72$ control and $n_2 = 4$ treated units. There are $T = 16$ observations per unit per outcome, $T_1 = 10$ of which are in the pre-intervention period. The objectives of the analysis are threefold. Firstly, we are interested in assessing the evidence for the existence of common factors underlying the four outcomes. Secondly, we are aiming to investigate the impact of the intervention on each of the four treated units (Derby, Enfield, Kingston upon Thames and Southwark) individually. Thirdly, we wish to assess the evidence for a non zero average intervention effect ϑ_{tk} ($t > 10$).

In order to achieve these goals, we fit our proposed MVFA+AR model. We set $k_1 = 20$ and $k_2 = 10$. We run MCMC for 1500000 iterations³; the first 250000 samples are discarded as a burn-in and we apply a thinning factor of 500 to the remaining draws. Therefore our MCMC sample consists of 2500 draws from the joint posterior of the model parameters. Convergence is assessed by visual inspection of posterior traceplots for some randomly chosen shrinkage parameters δ_{jk} , variance parameters ψ_{ik}^2 and the counterfactuals $y_{itk}^{(0)}$ ($i > 72$ and $t > 10$). These indicate that the chain has reached its stationary distribution. Further, we run an additional 9 chains and compare the posterior densities of these parameters with the ones obtained from the first chain. No major differences are found. Therefore, we conclude that the chain has converged.

For each outcome, we also fit the univariate FA model with the MGPS prior and $k_1 = 20$. However, the conclusions that we reached regarding the impact of the intervention are very similar (except for minor losses in precision of the causal effect estimates) to the ones obtained from the MVFA+AR model. Thus, the results from the FA model are not discussed further.

5.2 Results

Figure 6 summarises the posterior distribution of $\sum_{i=1}^{n_1+n_2} |\gamma_{ikj}|$ i.e. the L_1 -norm of the j -th column of the loadings matrix, where $\gamma_{ikj} = \lambda_{ik,j-k_1}$ for $j > k_1$. We see that the norm quickly decreases with j for both outcome-specific and shared factors, demonstrating the shrinkage effects of the MGPS prior. Inference on the number of non-negligible factors can be done as described in Section 3.4. If we use $m = 0.1$, then the median posterior number of non-negligible shared factors is 2 (95% CI [2,4]) and the

³This took approximately 9 hours on a Linux machine with an Intel i7-6700 3.4GHz CPU.

Results for time point $t = T_1 + 1$ and outcome $k = 1$ (of $K = 3$ outcomes)

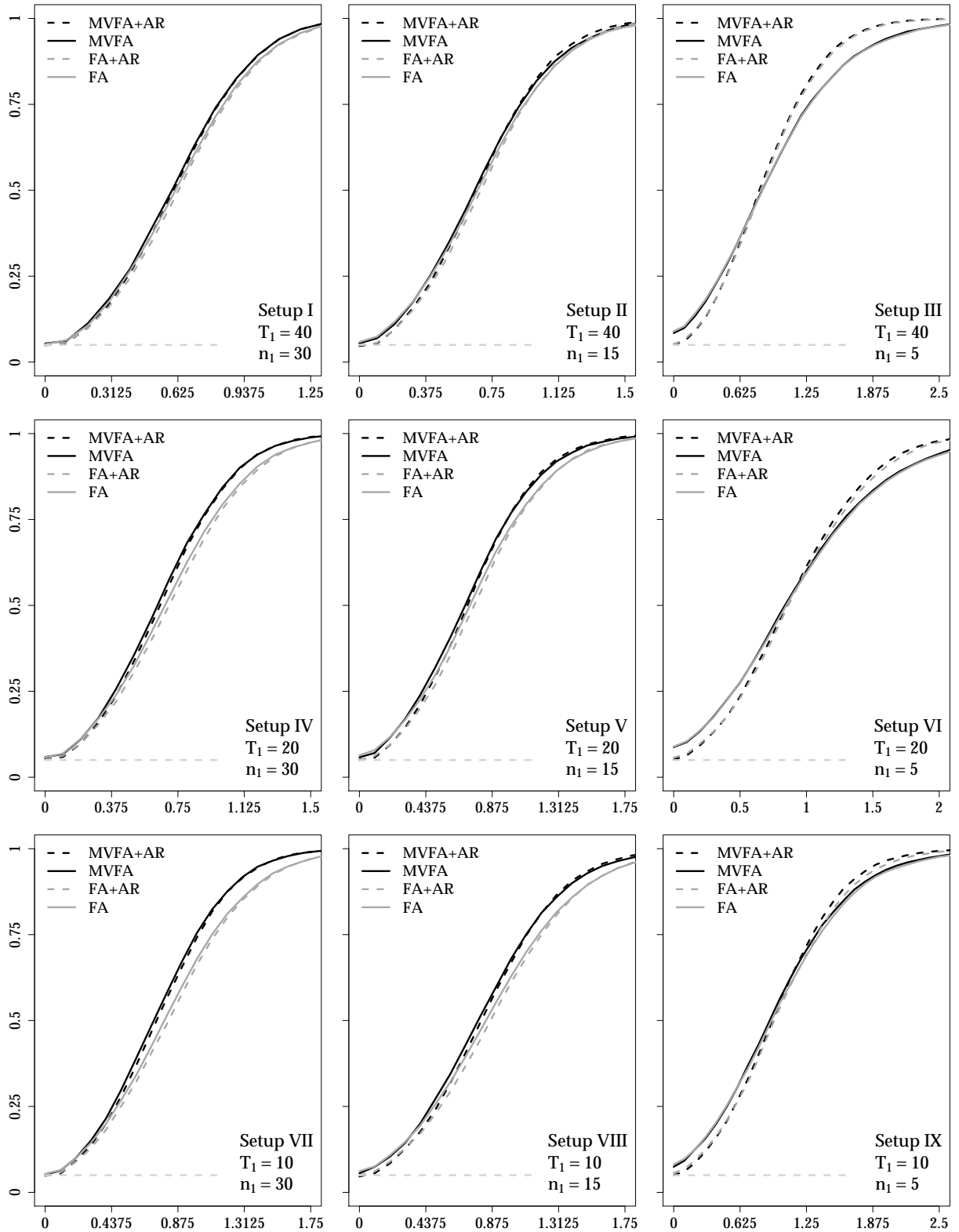


Figure 2: Results of the simulation study for the *first outcome* $k = 1$ and *first post-intervention time point* $t = T_1 + 1$. The figure presents the probability of detecting an intervention effect (y-axis) as a function of $\vartheta_{T_1+1,1}$ (x-axis) in Setups I-IX. All results are based on 10000 datasets simulated from the MVFA+AR model. The horizontal dashed lines indicate 5%, the desired detection rate when the intervention $\vartheta_{T_1+1,1} = 0$.

Results for time point $t = T$ and outcome $k = 1$ (of $K = 3$ outcomes)

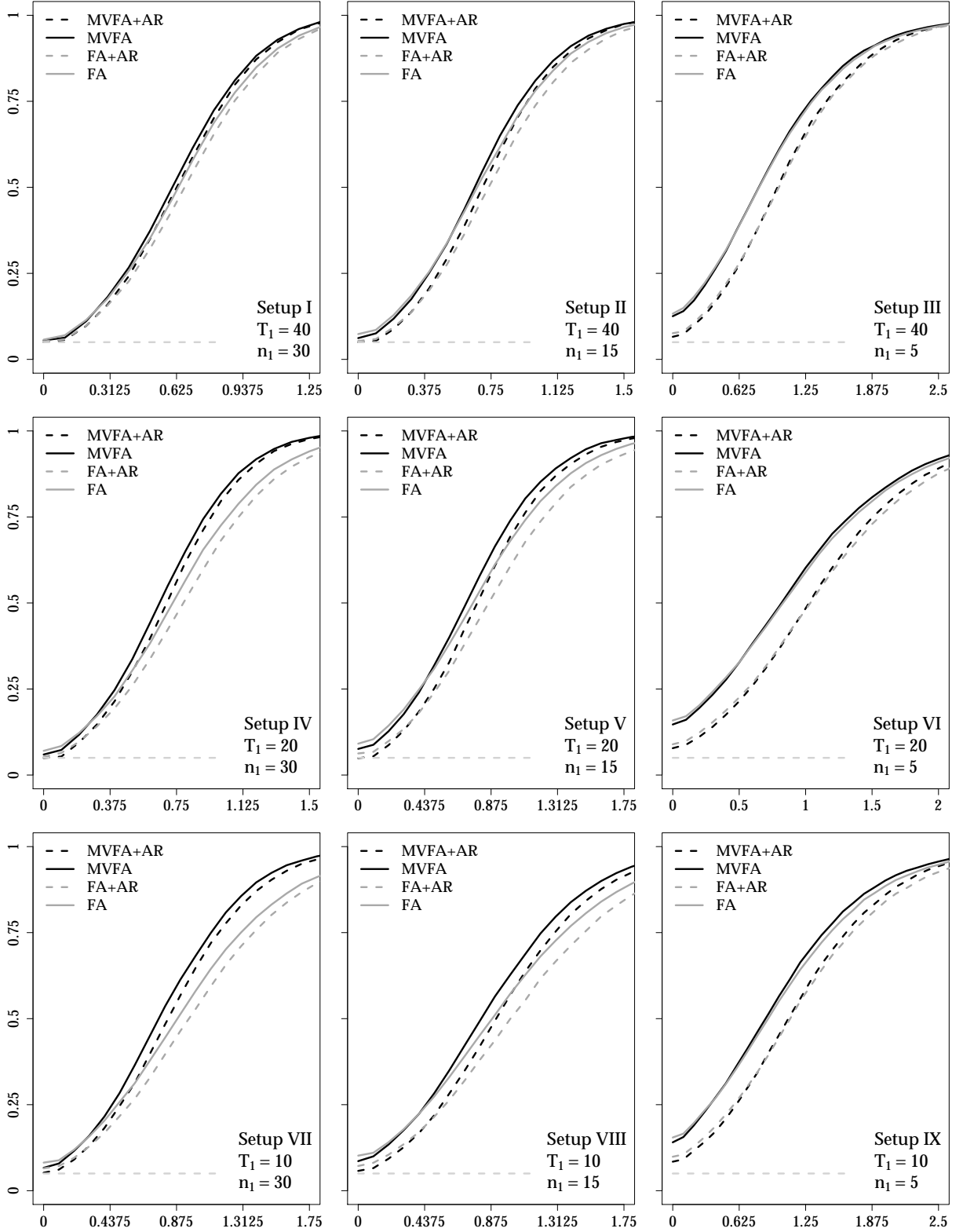


Figure 3: Results of the simulation study for the *first outcome* $k = 1$ and *last post-intervention time point* $t = T$. The figure presents the probability of detecting an intervention effect (y-axis) as a function of $\vartheta_{T,1}$ (x-axis) in Setups I-IX. All results are based on 10000 datasets simulated from the MVFA+AR model. The horizontal dashed lines indicate 5%, the desired detection rate when the intervention $\vartheta_{T,1} = 0$.

median posterior number of factors specific to outcomes 1-4 is 6 (95% CI [4,7]), 4 (95% CI [3,5]), 14 (95% CI [10,19]) and 11 (95% CI [9,14]), respectively.

There is not enough evidence in the data to support a significant intervention effect in each unit individually. This is evident in Figures 4 and 5 which show estimated counterfactuals along with their 95% credible intervals for Derby-Enfield and Kingston-Southwark, respectively. We see that for all treated units and outcomes, the estimated counterfactuals do not appear to be systematically higher than observed values. Further, the 95% credible intervals of the counterfactuals contain the observed values y_{itk} for the most of the combinations of i ($i > 72$), t ($t > 10$) and k . This suggests that the data are compatible with what would have been observed had intervention not taken place.

The point estimates of ϑ_{tk} , the average (over units) intervention effect, for all post-intervention time points 11-16 and outcomes, along with their 95% credible intervals, are shown in Table 2. We see that the credible intervals for antisocial behaviour, violent crimes and sexual crimes are nearly symmetrical about zero. Therefore we conclude that there is no evidence for an impact of the intervention on these outcomes. For hospital admissions, the point estimates are all negative (a negative values means admissions would be higher with no intervention). One of the advantages of the Bayesian approach is that it allows us to estimate many interesting causal quantities directly from the posterior distribution of the counterfactuals. Here we focus on the probability that $\vartheta_{tk} > 0$, which for hospital admissions and time points 11-16 is 0.41, 0.18, 0.16, 0.04, 0.10 and 0.07, respectively. Some of these values are low suggesting that the intervention succeeded in reducing the rate of hospital admissions. However, most of them are higher than 5% and therefore the evidence is inconclusive.

Setup	I	II	III	IV	V	VI	VII	VIII	IX
T_1	40	40	40	20	20	20	10	10	10
n_1	30	15	5	30	15	5	30	15	5

Results for $k = 1$ and $t = T_1 + 1$									
Bias									
MVFA+AR	0.017	0.030	0.119	0.034	0.056	0.155	0.072	0.103	0.195
MVFA	0.035	0.075	0.327	0.055	0.105	0.332	0.100	0.158	0.357
FA+AR	0.026	0.043	0.126	0.065	0.090	0.166	0.126	0.146	0.209
FA	0.044	0.087	0.324	0.082	0.131	0.330	0.144	0.186	0.351
Standard error									
MVFA+AR	0.314	0.354	0.480	0.347	0.388	0.508	0.396	0.440	0.539
MVFA	0.322	0.383	0.684	0.355	0.418	0.669	0.406	0.470	0.666
FA+AR	0.330	0.371	0.492	0.387	0.424	0.526	0.460	0.494	0.569
FA	0.335	0.398	0.691	0.392	0.449	0.677	0.466	0.517	0.680
Mean credible interval width									
MVFA+AR	1.245	1.394	1.904	1.370	1.548	2.055	1.628	1.833	2.286
MVFA	1.262	1.471	2.429	1.386	1.618	2.450	1.653	1.912	2.614
FA+AR	1.309	1.465	1.934	1.539	1.707	2.104	1.902	2.045	2.357
FA	1.317	1.517	2.425	1.532	1.736	2.462	1.898	2.076	2.634
False positive rate									
MVFA+AR	0.050	0.047	0.052	0.053	0.051	0.051	0.048	0.046	0.050
MVFA	0.054	0.055	0.085	0.059	0.058	0.088	0.051	0.056	0.074
FA+AR	0.048	0.049	0.052	0.051	0.049	0.057	0.049	0.050	0.057
FA	0.052	0.059	0.090	0.058	0.064	0.089	0.054	0.062	0.081

Results for $k = 1$ and $t = T$									
Bias									
MVFA+AR	0.014	0.032	0.208	0.043	0.074	0.263	0.113	0.170	0.333
MVFA	0.035	0.084	0.370	0.069	0.133	0.409	0.152	0.241	0.481
FA+AR	0.033	0.057	0.217	0.100	0.140	0.282	0.225	0.259	0.357
FA	0.054	0.107	0.372	0.121	0.186	0.420	0.246	0.308	0.490
Standard error									
MVFA+AR	0.326	0.375	0.676	0.374	0.437	0.737	0.484	0.569	0.800
MVFA	0.331	0.398	0.780	0.383	0.467	0.825	0.496	0.606	0.877
FA+AR	0.355	0.414	0.703	0.465	0.538	0.784	0.640	0.703	0.861
FA	0.358	0.433	0.800	0.469	0.554	0.858	0.642	0.717	0.913
Mean credible interval width									
MVFA+AR	1.284	1.485	2.483	1.485	1.746	2.660	1.913	2.216	2.931
MVFA	1.283	1.498	2.456	1.459	1.714	2.533	1.874	2.163	2.789
FA+AR	1.392	1.607	2.518	1.812	2.052	2.722	2.454	2.626	3.023
FA	1.373	1.574	2.462	1.728	1.920	2.570	2.341	2.471	2.863
False positive rate									
MVFA+AR	0.050	0.050	0.066	0.049	0.049	0.078	0.052	0.058	0.084
MVFA	0.055	0.062	0.125	0.060	0.076	0.147	0.067	0.086	0.141
FA+AR	0.052	0.054	0.076	0.053	0.063	0.089	0.065	0.072	0.099
FA	0.057	0.074	0.133	0.070	0.091	0.159	0.082	0.102	0.155

Table 1: Results of the simulation study for the *first outcome* $k = 1$ and post-intervention time points $t = T_1 + 1$ and $t = T$. The table presents the bias of the point estimates of ϑ_{t1} , the standard error of the point estimates, the mean width of the 95% credible intervals and the false positive rate. All results are based on 10000 simulated datasets from the MVFA+AR model.

Estimates of ϑ_{tk}				
t	Antisocial behaviour	Violent crimes	Hospital admissions	Sexual crimes
11	-0.09 [-0.26,0.07]	-0.1 [-0.29,0.09]	-1.2 [-13.1,9.8]	-0.003 [-0.013,0.007]
12	-0.04 [-0.27,0.19]	-0.12 [-0.44,0.21]	-5.6 [-18.6,6.7]	-0.003 [-0.018,0.011]
13	-0.04 [-0.36,0.28]	-0.22 [-0.63,0.23]	-6.9 [-20.9,7.2]	-0.008 [-0.026,0.011]
14	0.11 [-0.32,0.52]	-0.22 [-0.7,0.27]	-12.9 [-27.9,1.2]	-0.008 [-0.028,0.012]
15	0.26 [-0.19,0.73]	-0.09 [-0.6,0.45]	-9.8 [-25.1,5]	-0.01 [-0.031,0.012]
16	0.36 [-0.12,0.86]	-0.06 [-0.58,0.49]	-11.2 [-26.6,3.2]	-0.009 [-0.029,0.013]

Mean observed values				
t	Antisocial behaviour	Violent crimes	Hospital admissions	Sexual crimes
11	2.50	5.73	159.4	0.170
12	2.50	5.66	150.8	0.169
13	2.49	5.47	147.3	0.158
14	2.53	5.35	144.4	0.149
15	2.56	5.35	145.3	0.147
16	2.58	5.26	142.2	0.148

Table 2: Top panel: estimated average (over units) treatment effect for each outcome and post-intervention time point. The 95% posterior credible intervals are shown in brackets. Bottom panel: average (over units) observed values for each outcome and post-intervention time point.

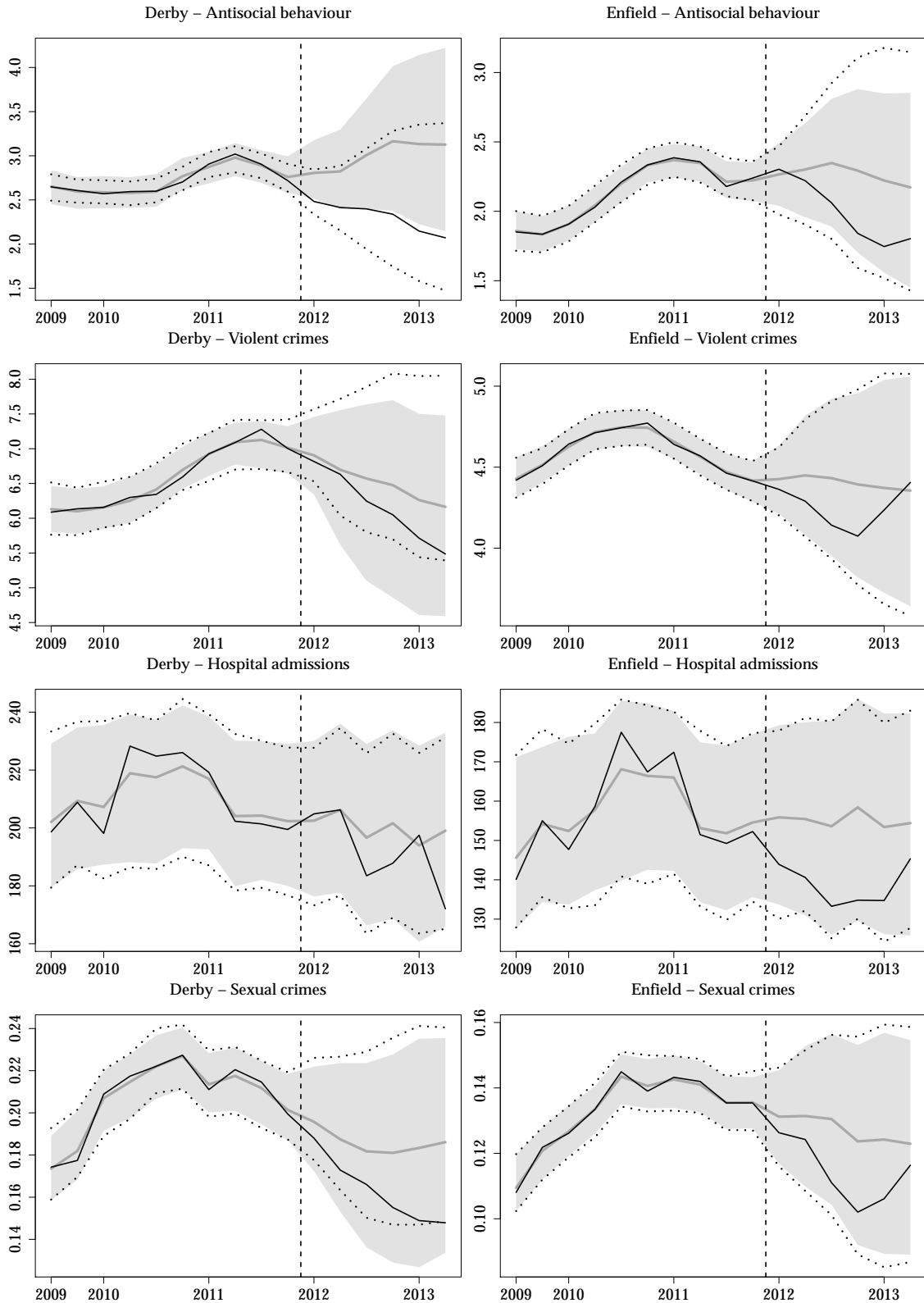


Figure 4: Results of the real data analysis for Derby and Enfield. The solid black lines represent the observed data. The solid grey lines indicate the posterior mean of $y_{itk}^{(0)}$ obtained by fitting the MVFA+AR model. The shaded grey regions represent the 95% credible intervals of $y_{itk}^{(0)}$ obtained from the same model. Finally, the dashed black lines indicate the 95% credible intervals for $y_{itk}^{(0)}$ obtained by analysing each outcome in turn with the FA model.

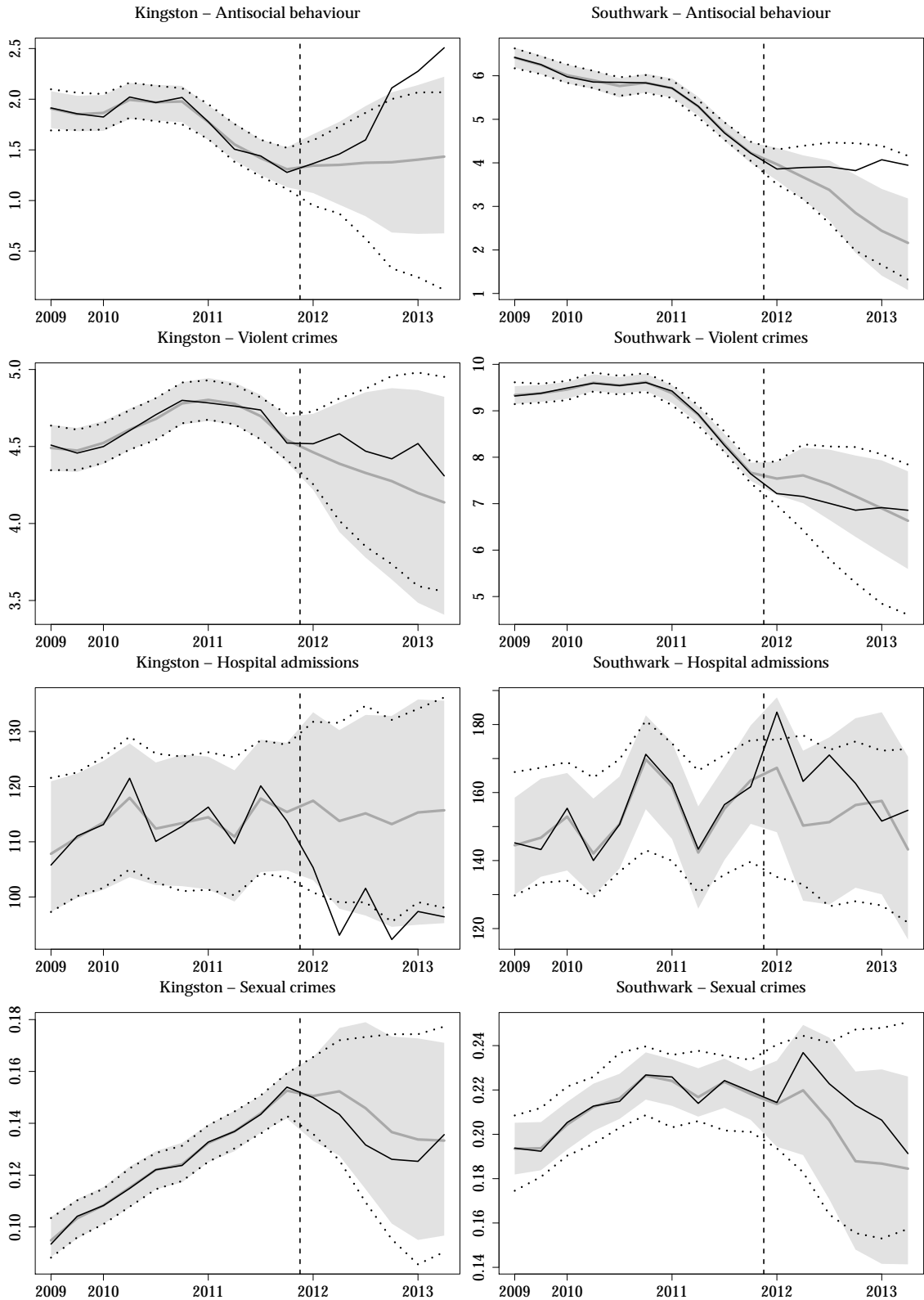


Figure 5: Results of the real data analysis for Kingston and Southwark. The solid black lines represent the observed data. The solid grey lines indicate the posterior mean of $y_{itk}^{(0)}$ obtained by fitting the MVFA+AR model. The shaded grey regions represent the 95% credible intervals of $y_{itk}^{(0)}$ obtained from the same model. Finally, the dashed black lines indicate the 95% credible intervals for $y_{itk}^{(0)}$ obtained by analysing each outcome in turn with the FA model.

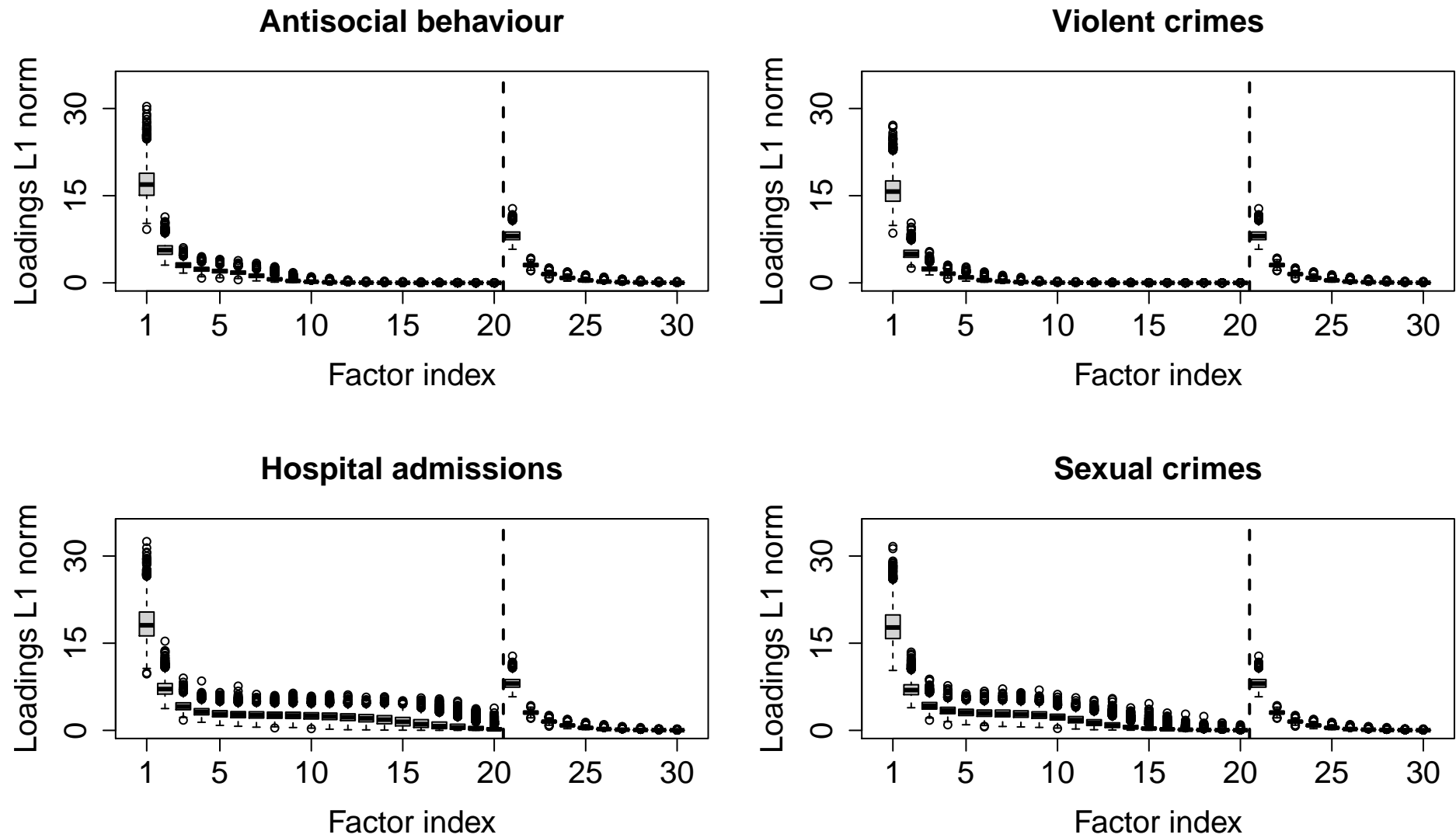


Figure 6: Analysis of the alcohol licensing dataset using the proposed MVFA model. The figure presents posterior boxplots of $\sum_{i=1}^n |\gamma_{ikj}|$, the L_1 -norm of the j -th column of the loadings matrix. The boxplots are based on an MCMC sample of size 2500. Factors 1-20 are outcome-specific whereas factors 21-30 are shared across outcomes.

6 Discussion

In this work, we have introduced the MVFA+AR model for evaluating the impact of a dichotomous intervention from time-series observational data. Our model extends in two ways the FA model frequently used for causal inference in this setting. First, it models multiple correlated outcomes jointly. Second, it accounts for autocorrelation within each of the outcomes. Both of these extensions enable more efficient estimation of the effect of an intervention on all, or any one, of the multiple outcomes. An important facet of the proposed model is that it provides posterior credible intervals for the causal effects of interest that account for the uncertainty about the number of factors.

The ability to make inference is inherent to the Bayesian framework and therefore to our method. This gives it an advantage over many existing approaches for causal inference using time-series observational data, including frequentist approaches based on the FA model (Gobillon and Magnac, 2016; Chan and Kwok, 2016; Xu, 2017; Li, 2018) and synthetic control-type approaches (Abadie *et al.*, 2010; Hsiao *et al.*, 2012; Doudchenko and Imbens, 2016; Ben-Michael *et al.*, 2018). The reason is that, to allow for inference, these methods require assumptions to allow for inference that might be unlikely to hold in some applications and therefore may yield confidence intervals that do not reflect the true uncertainty in the estimates of the causal effect. For example, the parametric bootstrap approach of Xu relies on the assumption that the error terms in the FA model are homoskedastic at each time point. For the approaches of Abadie *et al.* (2010) and Hsiao *et al.* (2012), inference is typically done with a “placebo test”, a procedure that is akin to a permutation test. However, the validity of this test is debatable unless one is willing to assume that the unit that received the intervention was chosen at random (Ben-Michael *et al.*, 2018). These assumptions are not essential for our method, suggesting it can be a useful alternative in applications where they are unlikely to hold.

Our simulation studies indicate that the estimates of intervention effects obtained from the MVFA+AR model are more precise than those obtained from the standard FA model. This can lead to considerable gains in power for detecting an intervention effect. Further, we found that the MVFA+AR model has better type I error rate control compared to the standard FA model. Both these gains occur when either the total number of pre-intervention time points or the total number of control units is relatively small. This is the case in many practical problems.

We applied our methodology to estimate the impact of CIPs on alcohol related harms. We found evidence for the existence of common factors driving the outcomes that we considered. We identified no major impact of CIPs on the rate of antisocial behaviour incidents, violent crimes and sexual crimes. The analysis provides some evidence that the intervention has led to a decrease in the rate of alcohol-related hospital admissions. However, the effect is not significant, i.e. the 95% credible intervals contain zero.

There are limitations to the proposed method. Our model allows for loadings that are shared across all outcomes. However, with $K > 2$ outcomes there is the possibility that there are loadings which are shared only between k of them, where $2 \leq k \leq K - 1$. There may be a benefit in extending the model to allow for loadings that are common only to a subset of the outcomes. Another extension that may be useful would replace the AR(1) structure that we assume with a more general ARMA process. However, this may not be feasible, given the length of the time-series encountered in many practical applications.

Several possible directions for future research exist. The proposed model does not make use of the geographical location of the units. Such information may be of value, since we expect the outcomes of units with spatial proximity to be correlated. It may be worth extending the model to account for this. Lopes *et al.* (2008) achieve this by assuming a spatial model for the loadings. Finally, although our model should perform well when the normality assumption holds approximately, it cannot be used when the data drastically deviate from this assumption, e.g. when the outcomes are dichotomous. Therefore, it is important to develop a model for mixed outcomes. We shall consider such extensions in our future research.

Acknowledgements

The authors thank Frank de Vocht for useful discussions regarding the real data application. This work is funded by the NIHR Health Protection Unit on Evaluation of Interventions (PS); Medical Research Council grants MC_UU_00002/10 (SRS) and MC_UU_00002/11 (DDA); and by Public Health England (DDA). This study is further funded by the National Institute for Health Research (NIHR) Programme Grants for Applied Research programme (Grant Reference Number RP-PG-0616-20008). The views expressed are those of the author(s) and not necessarily those of the NIHR or the Department of Health and Social Care

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**(490), 493–505.
- Angrist, J. D. and Pischke, J.-S. (2009). *Mostly harmless econometrics: An empiricist’s companion*. Princeton University Press.
- Avalos-Pacheco, A., Rossell, D., and Savage, R. S. (2018). Heterogeneous large datasets integration using Bayesian factor regression. *arXiv preprint arXiv:1810.09894*.
- Bai, J. (2009). Panel data models with interactive fixed effects. *Econometrica*, **77**(4), 1229–1279.
- Bai, J. and Ng, S. (2002). Determining the number of factors in approximate factor models. *Econometrica*, **70**(1), 191–221.
- Bartolucci, F., Pennoni, F., and Vittadini, G. (2016). Causal latent markov model for the comparison of multiple treatments in observational longitudinal studies. *Journal of Educational and Behavioral Statistics*, **41**(2), 146–179.
- Ben-Michael, E., Feller, A., and Rothstein, J. (2018). The augmented synthetic control method. *arXiv preprint arXiv:1811.04170*.
- Bhattacharya, A. and Dunson, D. B. (2011). Sparse bayesian infinite factor models. *Biometrika*, **98**(2), 291–306.
- Brodersen, K. H., Gallusser, F., Koehler, J., Remy, N., and Scott, S. L. (2015). Inferring causal impact using bayesian structural time-series models. *The Annals of Applied Statistics*, **9**(1), 247–274.
- Card, D. (1990). The impact of the mariel boatlift on the miami labor market. *Industrial and Labor Relation*, **43**(2), 245–257.
- Carvalho, C. M., Chang, J., Lucas, J. E., Nevins, J. R., Wang, Q., and West, M. (2008). High-dimensional sparse factor modeling: applications in gene expression genomics. *Journal of the American Statistical Association*, **103**(484), 1438–1456.
- Cavallo, E., Galiani, S., Noy, I., and Pantano, J. (2013). Catastrophic natural disasters and economic growth. *The Review of Economics and Statistics*, **95**(5), 1549–1561.
- Chan, M. and Kwok, S. (2016). Policy evaluation with interactive fixed effects.
- De Vito, R., Bellio, R., Trippa, L., and Parmigiani, G. (2016). Multi-study factor analysis. *arXiv preprint arXiv:1611.06350*.
- De Vito, R., Bellio, R., Trippa, L., and Parmigiani, G. (2018). Bayesian multi-study factor analysis for high-throughput biological data. *arXiv preprint arXiv:1806.09896*.
- de Vocht, F. (2016). Inferring the 1985–2014 impact of mobile phone use on selected brain cancer subtypes using Bayesian structural time series and synthetic controls. *Environment International*, **97**, 100–107.
- de Vocht, F., Tilling, K., Pliakas, T., Angus, C., Egan, M., Brennan, A., Campbell, R., and Hickman, M. (2017). The intervention effect of local alcohol licensing policies on hospital admission and crime: a natural experiment using a novel Bayesian synthetic time-series method. *Journal of Epidemiology & Community Health*, **71**(9), 912–918.
- Doudchenko, N. and Imbens, G. W. (2016). Balancing, Regression, Difference-In-Differences and Synthetic Control Methods: A Synthesis. *arXiv preprint arXiv:1610.07748*.
- Geweke, J. and Zhou, G. (1996). Measuring the pricing error of the arbitrage pricing theory. *The Review of Financial Studies*, **9**(2), 557–587.
- Gobillon, L. and Magnac, T. (2016). Regional policy evaluation: Interactive fixed effects and synthetic controls. *Review of Economics and Statistics*, **98**(3), 535–551.

- Holland, P. W. (1986). Statistics and causal inference. *Journal of the American statistical Association*, **81**(396), 945–960.
- Hsiao, C., Ching, S. H., and Wan, S. K. (2012). A panel data approach for program evaluation: measuring the benefits of political and economic integration of Hong Kong with mainland China. *Journal of Applied Econometrics*, **27**(5), 705–740.
- Jones, A. M. and Rice, N. (2011). Econometric evaluation of health policies. In S. Glied and P. Smith, editors, *The Oxford Handbook of Health Economics*, Oxford Handbooks, pages 890–923. OUP Oxford.
- Kastner, G. and Frühwirth-Schnatter, S. (2014). Ancillarity-sufficiency interweaving strategy (ASIS) for boosting MCMC estimation of stochastic volatility models. *Computational Statistics and Data Analysis*, **76**(C), 408–423.
- Lanza, S. T., Coffman, D. L., and Xu, S. (2013). Causal inference in latent class analysis. *Structural Equation Modeling: a Multidisciplinary Journal*, **20**(3), 361–383.
- Li, K. (2018). Inference for factor model based average treatment effects. Available at SSRN 3112775.
- Lopes, H. F. and West, M. (2004). Bayesian model assessment in factor analysis. *Statistica Sinica*, **14**, 41–67.
- Lopes, H. F., Salazar, E., Gamerman, D., et al. (2008). Spatial dynamic factor analysis. *Bayesian Analysis*, **3**(4), 759–792.
- McAlinn, K., Aastveit, K. A., Nakajima, J., and West, M. (2017). Multivariate Bayesian predictive synthesis in macroeconomic forecasting. *arXiv preprint arXiv:1711.01667*.
- Montagna, S., Tokdar, S. T., Neelon, B., and Dunson, D. B. (2012). Bayesian latent factor regression for functional and longitudinal data. *Biometrics*, **68**(4), 1064–1073.
- Montagna, S., Irincheeva, I., and Tokdar, S. T. (2018a). High-dimensional bayesian fourier analysis for detecting circadian gene expressions. *arXiv preprint arXiv:1809.04347*.
- Montagna, S., Wager, T., Barrett, L. F., Johnson, T. D., and Nichols, T. E. (2018b). Spatial bayesian latent factor regression modeling of coordinate-based meta-analysis data. *Biometrics*, **74**(1), 342–353.
- Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective*. MIT Press.
- Robbins, M. W., Saunders, J., and Kilmer, B. (2017). A framework for synthetic control methods with high-dimensional, micro-level data: Evaluating a neighbourhood-specific crime intervention. *Journal of the American Statistical Association*, **112**(517), 109–126.
- Robins, J. M., Hernan, M. A., and Brumback, B. (2000). Marginal structural models and causal inference in epidemiology. *Epidemiology*, **11**(5).
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, **66**(5), 688.
- Rue, H. (2001). Fast sampling of Gaussian Markov random fields. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **63**(2), 325–338.
- Samartsidis, P., Seaman, S. R., Presanis, A. M., Hickman, M., and De Angelis, D. (2019). Assessing the causal effect of binary interventions from observational panel data with few treated units. *Statistical Science*.
- Srivastava, S., Engelhardt, B. E., and Dunson, D. B. (2017). Expandable factor analysis. *Biometrika*, **104**(3), 649–663.
- Tullio, F. and Bartolucci, F. (2019). Evaluating time-varying treatment effects in latent markov models: An application to the effect of remittances on poverty dynamics.
- Xu, Y. (2017). Generalized synthetic control method: Causal inference with interactive fixed effects models. *Political Analysis*, **25**(1), 5776.