

New Rover

*hopefully written
letter for you B*

Here now we deal with understanding the drpm package.
It seems that there is just one function for fitting a
model: drpm-fit.
The problem is understanding the meaning of the
input parameters, from the user.

y	An m x T numeric matrix containing the response values measured over time.
s_coords	Two-column matrix containing spatial locations (i.e., longitude and latitude).
M	Parameter related to Dirichlet process scale or dispersion parameter.
initial_partition	The initial partition that is elicited from an expert.
starting_alpha	Initial value for alpha. If alpha_0 = TRUE (see below), then alpha is fixed to this value.
unit_specific_alpha	Logical that indicates if a unique alpha is to be included for each unit.
time_specific_alpha	Logical that indicates if a unique alpha is to be included for each time point.
alpha_0	Logical where TRUE indicates alpha_t = starting_alpha (i.e., alpha is not updated)
eta1_0	Logical where TRUE indicates eta1 = 0 (i.e., observations are conditionally independent).
phi1_0	Logical where TRUE indicates phi1 = 0 (i.e., iid model for atoms)
modelPriors	Vector containing the following hierarchical model prior values in order

```
drpm_fit(y,
          s_coords=NULL,
          M=1,
          initial_partition = NULL,
          starting_alpha = 0.5,
          unit_specific_alpha = FALSE,
          time_specific_alpha = FALSE,
          alpha_0=FALSE,
          eta1_0=FALSE,
          phi1_0=FALSE,
          modelPriors=c(0,100^2,1,1,1,1),
          alphaPriors=rbind(c(1,1)),
          simpleModel = 0,
          theta_tau2 =c(0, 2),
          SpatialCohesion=4,
          cParms=c(0, 1, 2, 1),
          mh=c(0.5, 1, 0.1, 0.1, 0.1),
          verbose=FALSE,
          draws=1100,burn=100,thin=1)
```

- m0 - mean of phi0
- s20 - variance of phi0
- A - upper bound on sigma
- At - upper bound on tau
- Al - upper bound on lambda
- be - scale prior for xi

next page for better explanation &

alphaPriors	Matrix containing beta prior values for the alpha parameters. The columns correspond to the 1st and 2nd shape beta shape parameters. The number of rows corresponds to the number of units. If no initial partition is supplied then only the first row is used.
simpleModel	If 1 then the simple model in Sally's simulation is employed.
theta_tau2	Only used if simpleModel=1. This is a two-dimensional vector with the first entry corresponding to the value that is fixed to theta and the second to tau2
SpatialCohesion	Scalar indicating what cohesion to use if spatial coordinates are provided <ul style="list-style-type: none"> 3 - auxiliary similarity 4 - double dipper similarity
cParms	Vector containing spatial cohesion parameter values (see below for more details) <ul style="list-style-type: none"> mu0 - center of NNIG k0 - scale parameter of gaussian part of NNIG v0 - degrees of freedom for inverse-gamma part of NNIG L0 - scale parameter for inverse-gamma part of NNIG
mh	Vector of gaussian standard deviations for metropolis updates of sigma2, tau, lambda, eta1, phi1.
draws	Number of MCMC samples to collect
burn	Number of the MCMC samples discarded in the burn-in phase of the sampler
thin	The amount of thinning desired for the chain
verbose	Logical indicating if MCMC progression should be print to screen

"M=4" does not mean "we want 4 clusters one", the meaning/reason seems more complex
 M turns how well max clusters will fit the model tend to build

y → An m x T numeric matrix containing the response values measured over time.
 s_coords → Two-column matrix containing spatial locations (i.e., longitude and latitude).
 M → Parameter related to Dirichlet process scale or dispersion parameter.
 initial_partition → caved be set to null if there is no information
 The initial partition that is elicited from an expert.

starting_alpha → Initial value for alpha. If alpha_0 = TRUE (see below), then alpha is fixed to this value.
 unit_specific_alpha → Logical that indicates if a unique alpha is to be included for each unit. *if FALSE => each unit gets its own or => one*
 time_specific_alpha → Logical that indicates if a unique alpha is to be included for each time point. *if FALSE => at each time we get different or*
 alpha_0 → Logical where TRUE indicates alpha_t = starting_alpha (i.e., alpha is not updated) *at*
 eta1_0 → Logical where TRUE indicates eta1 = 0 (i.e., observations are conditionally independent).
 phi1_0 → Logical where TRUE indicates phi1 = 0 (i.e., iid model for atoms)
 modelPriors → Vector containing the following hierarchical model prior values in order

- m0 - mean of phi0
- s20 - variance of phi0
- A - upper bound on sigma
- At - upper bound on tau
- Al - upper bound on lambda
- be - scale prior for xi

drpm_fit(y,
 s_coords=NULL,
 M=1,
 initial_partition = NULL,
 starting_alpha = 0.5,
 unit_specific_alpha = FALSE,
 time_specific_alpha = FALSE,
 alpha_0=FALSE,
 eta1_0=FALSE,
 phi1_0=FALSE,
 modelPriors=c(0,100^2,1,1,1,1),

Dr + 10
 Dr + 6
 Come monica sei credo

```

alphaPriors=rbind(c(1,1)),
simpleModel = 0,
theta_tau2 =c(0, 2),
SpatialCohesion=4,
cParms=c(0, 1, 2, 1),
mh=c(0.5, 1, 0.1, 0.1, 0.1),
verbose=FALSE,
draws=1100,burn=100,thin=1)

```

not yet clarified

alphaPriors
simpleModel
theta_tau2

Matrix containing beta prior values for the alpha parameters. The columns correspond to the 1st and 2nd shape beta shape parameters. The number of rows corresponds to the number of units. If no initial partition is supplied then only the first row is used.

If 1 then the simple model in Sally's simulation is employed.

Only used if simpleModel=1. This is a two-dimensional vector with the first entry corresponding to the value that is fixed to theta and the second to tau2

SpatialCohesion

Scalar indicating what cohesion to use if spatial coordinates are provided

- 3 - auxiliary similarity
- 4 - double dipper similarity

cParms

Vector containing spatial cohesion parameter values (see below for more details)

- mu0 - center of NNIG
- k0 - scale parameter of gaussian part of NNIG
- v0 degrees of freedom for inverse-gamma part of NNIG
- L0 - scale parameter for inverse-gamma part of NNIG

mh
draws
burn
thin
verbose

Vector of gaussian standard deviations for metropolis updates of sigma2, tau, lambda, eta1, phi1.

Number of MCMC samples to collect

Number of the MCMC samples discarded in the burn-in phase of the sampler

The amount of thinning desired for the chain

Logical indicating if MCMC progression should print to screen

$$C_3(S_h, s_h^*) = M \times \Gamma(|S_h|) \times \int \prod_{i \in S_h} q(s_i | \xi_h) q(\xi_h) d\xi_h. \quad (3.4)$$

In Bayesian modeling $\int \prod_{i \in S_h} q(s_i | \xi_h) q(\xi_h) d\xi_h$ is often called the marginal likelihood or prior predictive distribution and is used to measure the similarity among the locations belonging to cluster h . Therefore, C_3 favors partitioned location vectors (s^*) that produce large marginal likelihood values. To simplify evaluating C_3 and retain coherence across sample sizes, $q(s|\xi)$ and $q(\xi)$ are specified to form a conjugate probability model. We emphasize however that we are not assuming the s_i 's to be random, we are simply employing the conjugate model as a means to measure spatial proximity and encourage co-clustering of locations that are near each other. Both areal and point referenced data can be considered when C_3 is employed, all that is required is specifying appropriate $q(s|\xi)$ and $q(\xi)$. For example, if point referenced data are available, a conjugate Gaussian/Gaussian-Inverse-Wishart model would be appropriate. In this case $\xi = (\mathbf{m}, \mathbf{V})$ would denote a mean and covariance, $q(s|\xi) = N(s|\mathbf{m}, \mathbf{V})$ a bivariate Gaussian density and $q(\xi) = NIW(\mathbf{m}, \mathbf{V} | \mu_0, \kappa_0, \nu_0, \Lambda_0)$ a bivariate Normal-Inverse-Wishart density. For areal data a conjugate multinomial/Dirichlet model could be utilized. In what follows we focus on point reference case and will occasionally refer to C_3 as the auxiliary cohesion. Finally, as in the previous two cohesions, $M \times \Gamma(|S_h|)$ is included to avoid creating many singleton clusters.

The fourth and final cohesion that we consider is similar to what Quintana et al. (In press) call a double dipper cohesion. It has the same form as C_3 , but instead of employing a prior predictive conjugate model, a posterior predictive conjugate model is used. Therefore C_4 has the following form

$$C_4(S_h, s_h^*) = M \times \Gamma(|S_h|) \times \int \prod_{i \in S_h} q(s_i | \xi_h) q(\xi_h | s_h^*) d\xi_h. \quad (3.5)$$

Since the posterior predictive is typically more peaked than the prior predictive, C_4 puts more weight on partitions that are local. Once again both areal and point referenced data are possible, but in what follows we focus on point-referenced and use the following conjugate model: $N_2(s_i | \mathbf{m}_h, \mathbf{V}_h) NIW(\mathbf{m}_h, \mathbf{V}_h | s_h^*)$.

In order to introduce the sPPM, let s_i denote the spatial coordinates of the i th item (note that these coordinates do not change over time) and let s_{jt}^* be the subset of spatial coordinates that belong to the j th cluster at time t . Then we express the EPPF of the t th partition with the following product form

$$\Pr(\rho_t | v_0, M) \propto \prod_{j=1}^{k_t} c(S_{jt}|M) g(s_{jt}^* | v_0). \quad (11)$$

Here, $c(\cdot|M) \geq 0$ is called the cohesion and is a set function that produces cluster weights a priori. We consider the cohesion $c(S_{jt}|M) = M \times (|S_{jt}| - 1)!$ as it has connections with the CRP making this version of the sPPM a type of spatially reweighted CRP. The similarity function $g(\cdot|v_0)$ is a set function parametrized by v_0 that measures the compactness of the spatial coordinates in s_{jt}^* , producing higher values if the spatial coordinates in s_{jt}^* are close to each other. Not all similarity functions preserve sample size consistency so to ensure this, after standardizing spatial locations, we employ

$$g(s_{jt}^* | v_0) = \int \prod_{i \in S_{jt}} N(s_i | \mathbf{m}, \mathbf{V}) NIW(\mathbf{m}, \mathbf{V} | \mathbf{0}, 1, v_0, I) dm d\mathbf{V}, \quad (12)$$

$$c(S_{jt}|M) \propto c(S_{jt}) \cdot f(\sigma_{jt}^*)$$

cluster in cluster h

the cohesion chosen the
we need wt works good
we need us to response
the correlations
 $c(S_{jt}) = \dots$ see the
 $f(\sigma_{jt}^*) = \dots$ see the
unreal

see
user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40

see user 3
user 40



Dependent Modeling of Temporal Sequences of Random Partitions

Garrett L. Page^{a,b}, Fernando A. Quintana^c, and David B. Dahl^a

^aDepartment of Statistics, Brigham Young University, Provo, UT; ^bBCAM—Basque Center of Applied Mathematics, Bilbao, Spain; ^cDepartamento de Estadística, Pontificia Universidad Católica de Chile, Millennium Nucleus Center for the Discovery of Structures in Complex Data, Santiago, Chile

ABSTRACT

We consider modeling a dependent sequence of random partitions. It is well known in Bayesian nonparametrics that a random measure of discrete type induces a distribution over random partitions. The community has therefore assumed that the best approach to obtain a dependent sequence of random partitions is through modeling dependent random measures. We argue that this approach is problematic and show that the random partition model induced by dependent Bayesian nonparametric priors exhibits counter-intuitive dependence among partitions even though the dependence for the sequence of random probability measures is intuitive. Because of this, we suggest directly modeling the sequence of random partitions when clustering is of principal interest. To this end, we develop a class of dependent random partition models that explicitly models dependence in a sequence of partitions. We derive conditional and marginal properties of the joint partition model and devise computational strategies when employing the method in Bayesian modeling. In the case of temporal dependence, we demonstrate through simulation how the methodology produces partitions that evolve gently and naturally over time. We further illustrate the utility of the method by applying it to an environmental dataset that exhibits spatio-temporal dependence. Supplemental files for this article are available online.

ARTICLE HISTORY

Received July 2020
Revised September 2021

KEYWORDS

Bayesian nonparametrics;
Correlated partitions;
Hierarchical Bayes modeling;
Spatio-temporal clustering

1. Introduction

We introduce a method to directly model dependence in a sequence of random partitions. Our approach is motivated by the practical problem of defining a prior distribution to model a sequence of random partitions that potentially exhibits substantial concordance over time (e.g., gently evolving clusterings over time). Traditionally, dependencies in random partitions (i.e., the clustering of units) have been obtained as a by-product of dependent random measures in Bayesian nonparametric (BNP) methods. We argue, however, that when a sequence of partitions is the inferential object of interest, then the sequence of partitions should be modeled *directly* rather than relying on *induced* random partition models, such as those implied by temporally dependent BNP models. But first, we review the literature on dependent BNP methods.

A nonexhaustive list of BNP methods that temporally correlate a sequence of random probability measures include Nieto-Barajas et al. (2012), Antoniano-Villalobos and Walker (2016), Gutiérrez, Mena, and Ruggiero (2016), Jo et al. (2017), Kalli and Griffin (2018), DeYoreo and Kottas (2018), and De Iorio et al. (2019). A common aspect of all these methods is that temporal dependence is accommodated in the sequence of random measures by way of the atoms or weights of the stick-breaking representation (Sethuraman 1994). An alternative approach to producing a sequence of temporally correlated random probability measures can be found in Caron, Davy, and Doucet (2007) and

Caron et al. (2017). Their construction is based on a generalised Pólya urn scheme where dependencies between distributions that evolve over time are induced by urn-like operations on counts and the parameters to which they are associated. A key insight associated with all mentioned approaches, however, is that the induced random partitions only exhibit weak dependence even when a sequence of random probability measures is highly correlated. To illustrate this point, we conducted a small Monte Carlo simulation where an induced sequence of partitions was generated with 10 time points and 20 units using the method of Caron et al. (2017). To measure similarity of partitions at different time points, we use a time-lagged Adjusted Rand Index (ARI) (Hubert and Arabie 1985). Figure 1 shows these values averaged over 10,000 Monte Carlo samples. Notice that as the temporal dependence parameter (α) increases, the partitions from time period t to $t + 1$ only become slightly more similar, such that the dependence between partitions is, at best, only weak. Further, the dependence is not temporally intuitive as it does not decay as a function of lagged time. In contrast, compare the dependence structure in Figure 1 to that in Figure 2, which contains average lagged ARI values between time-lagged partitions generated using the method developed in this article. Notice that unlike the induced random partitions generated using the method in Caron et al. (2017), the sequence of partitions generated using our approach displays intuitive temporal dependence. That is, as the time lag increases,

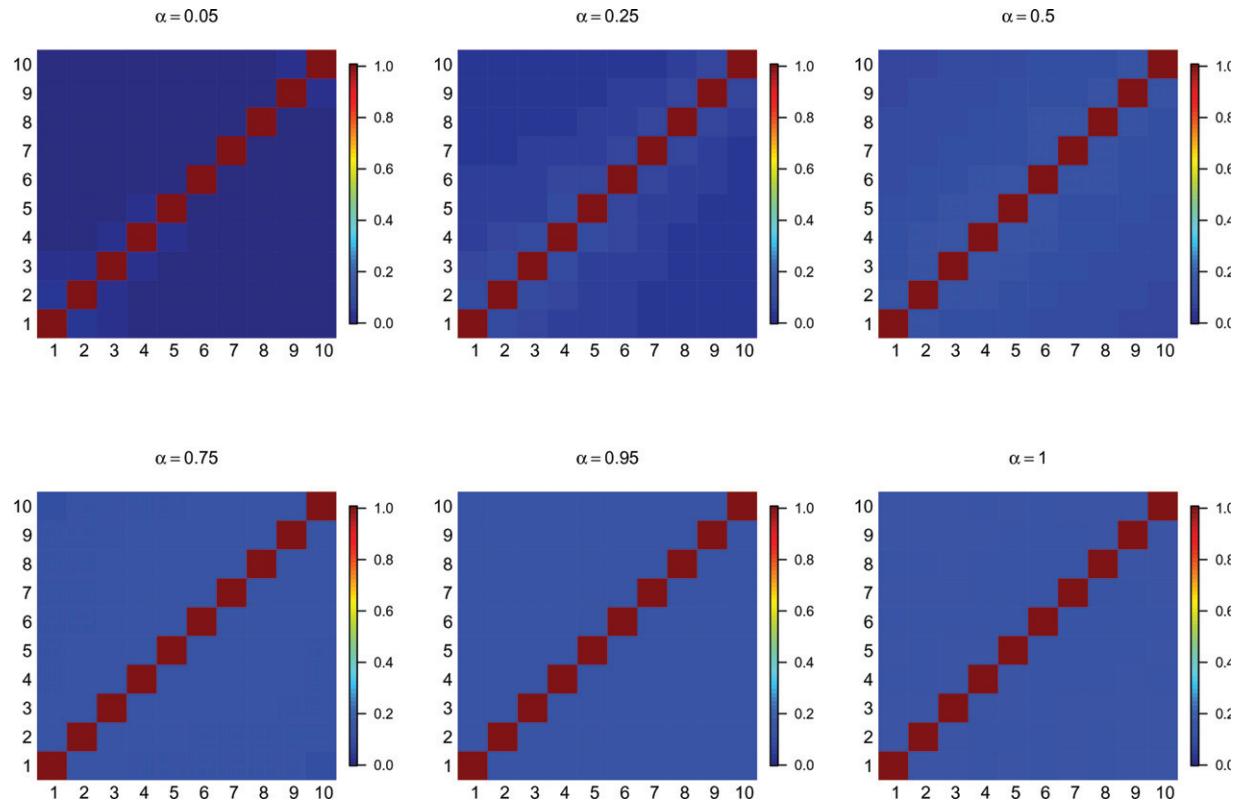


Figure 1. For various values of the temporal dependence parameter α , these plots show the lagged ARI values using the method of Caron et al. (2017) based on concentration parameter $M = 0.5$, discount parameter set to zero, and 10,000 Monte Carlo samples.

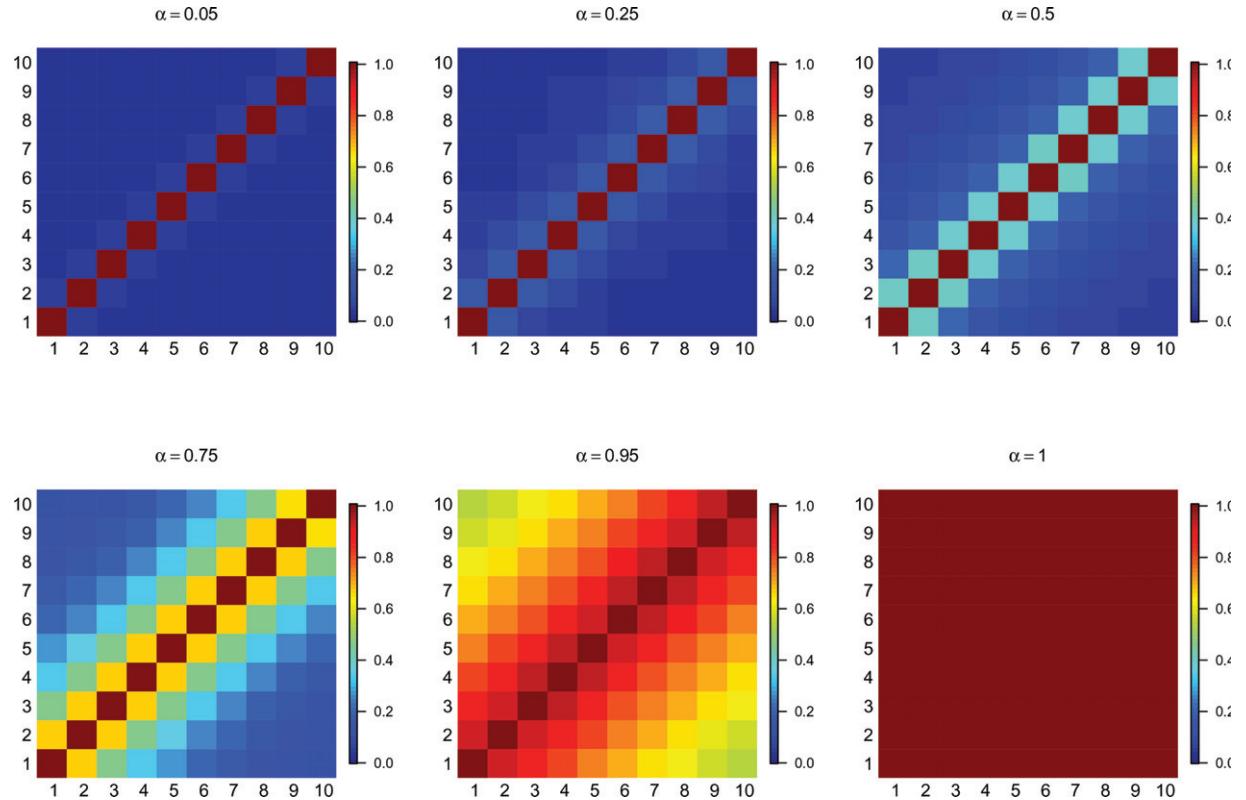


Figure 2. Lagged ARI values based on 10,000 Monte Carlo samples using the method developed in this paper. Partitions show natural and intuitive temporal dependence as lagged time increases and as the temporal dependence parameter α increases.

similarity between partitions decreases and as the temporal dependence parameter α increases, the partitions become dissimilar over time at a slower rate.

The counter-intuitive behavior displayed in Figure 1 is not unique to the approach of Caron et al. (2017). As noted by Wade, Walker, and Petrone (2014), the same type of behavior is

present when using a linear-dependent Dirichlet process mixture model. In fact, all BNP methods that model a sequence of discrete random probability measures will induce a random partition model with similarly weak correlation behavior. This behavior is analogous to trying to induce dependence among random variables from distributions with correlated parameters. There is no guarantee that correlated parameters would produce strong correlations among the random variables themselves.

Our approach is to consider the sequence of partitions as the parameter of principal interest and develop a method that models it directly. This will provide more control over how “smoothly” partitions evolve over time. Perhaps the work closest to ours, in the sense of explicitly modeling a sequence of partitions, can be found in Zanini et al. (2019). Their modeling approach for a temporally-referenced sequence of partitions can be applied to only two time points and differs from ours in that they do not focus on smooth evolution of partitions over time.

The remainder of the article is organized as follows. In Section 2 we present the proposed approach for a sequence of dependent random partitions, discuss its main properties, and suitable computational strategies for inference based on posterior simulation. Section 3 contains the results from three simulation studies that further explore aspects of the model. Section 4 describes an environmental data application and some concluding remarks are provided in Section 5. An accompanying Supplementary Materials file collects the proofs of results stated below, provides details on posterior simulation algorithms, and contains further simulation and data analysis results.

2. Joint Model for a Sequence of Partitions

Before detailing our method, we introduce some general notation. Let $i = 1, \dots, m$ denote the m experimental units at time t for $t = 1, \dots, T$. Let $\rho_t = \{S_{1t}, \dots, S_{kt}\}$ denote a partition of the m experimental units at time $t = 1, \dots, T$ into k_t clusters. An alternative notation is based on m cluster labels at time t denoted by $c_t = \{c_{1t}, \dots, c_{mt}\}$ where $c_{it} = j$ implies that $i \in S_{jt}$. Finally, any quantity with a “ $*$ ” superscript will be cluster-specific. For example, we will use μ_{jt}^* to denote the mean of cluster j at time t so that $\mu_{it} = \mu_{c_{it}t}^*$.

2.1. Temporal Modeling for Sequences of Partitions

Introducing temporal dependence in a collection of partitions requires formulating a joint probability model for (ρ_1, \dots, ρ_T) . Generically, we will denote this joint model with $\Pr(\rho_1, \dots, \rho_T)$. Temporal dependence among the ρ_t 's implies that the cluster configuration in ρ_t could be impacted by that found in $\rho_{t-1}, \rho_{t-2}, \dots, \rho_1$. However, we assume that the probability model for the sequence of partitions has a first-order Markovian structure. That is, the conditional distribution of ρ_t given $\rho_{t-1}, \rho_{t-2}, \dots, \rho_1$ only depends on ρ_{t-1} . Thus, we construct $\Pr(\rho_1, \dots, \rho_T)$ as

$$\Pr(\rho_1, \dots, \rho_T) = \Pr(\rho_T | \rho_{T-1}) \Pr(\rho_{T-1} | \rho_{T-2}) \cdots \Pr(\rho_2 | \rho_1) \Pr(\rho_1). \quad (1)$$

Here $\Pr(\rho_1)$ is an exchangeable partition probability function (EPPF) that describes how the m experimental units at time period 1 are grouped into k_1 distinct groups with frequencies n_{11}, \dots, n_{1k_1} . One characteristic of an EPPF that will prove useful in what follows is sample size consistency, or what De Blasi et al. (2015) refer to as the *addition rule*. This property dictates that marginalizing the last of $m + 1$ elements leads to the same model as if we only had m elements. A commonly encountered EPPF is that induced by a Dirichlet process (DP). This particular EPPF is sometimes referred to as a Chinese restaurant process (CRP) which can be seen as a special case from the family of product partition models (PPM). For more details, see De Blasi et al. (2015). Because we employ the EPPF of the CRP in what follows, we provide its form here

$$\Pr(\rho | M) = \frac{M^k}{\prod_{i=1}^n (M+i-1)} \prod_{i=1}^k (|S_i|-1)! \quad (2)$$

where k is the number of clusters in ρ and M is a concentration parameter controlling the number of clusters. We will denote this random partition distribution as $\text{CRP}(M)$.

Although conceptually straightforward, (1) is silent regarding how ρ_{t-1} influences the form of ρ_t . To make this explicit, we introduce an auxiliary variable that guides the similarity between ρ_t and ρ_{t-1} . Note that if two partitions are highly dependent, then the cluster configurations between them will change very little and as a result only a few of the m experimental units will change cluster assignment. Conversely, two partitions that exhibit low dependence will likely be comprised of very different cluster configurations. The auxiliary variable we introduce identifies which of the experimental units at time $t - 1$ will be considered for possible cluster reallocation at time t . Specifically, let γ_{it} denote the following

$$\gamma_{it} = \begin{cases} 1 & \text{if unit } i \text{ is not reallocated when moving from time } t-1 \text{ to } t \\ 0 & \text{otherwise,} \end{cases} \quad (3)$$

for $i = 1, \dots, m$. Notice that when $\gamma_{it} = 0$, item i is subject to reallocation at time t , but still may, by random assignment, end up in the same cluster at time t as it was at time $t - 1$. By construction, we set $\gamma_{i1} = 0$ for all i , that is, all experimental units are allocated to clusters during the first time period. We

then assume that $\gamma_{it} \stackrel{\text{ind}}{\sim} \text{Ber}(\alpha_t)$. Note that each of the $\alpha_t \in [0, 1]$ acts as a temporal dependence parameter. Specifically, we will interpret $\alpha_t = 1$ as implying that $\rho_t = \rho_{t-1}$ with probability 1. Conversely, when $\alpha_t = 0$, then ρ_t is independent of ρ_{t-1} . Further, when α_t is constant for all t , the degree of dependence among partitions is constant over time, whereas general values for α_t provide for varying degrees of dependence and more flexible partition patterns over time. For notational convenience, we introduce $\boldsymbol{\gamma}_t = (\gamma_{1t}, \gamma_{2t}, \dots, \gamma_{mt})$ which is an m -tuple comprised of zeros and ones. The augmented joint model changes (1) to

$$\Pr(\boldsymbol{\gamma}_1, \rho_1, \dots, \boldsymbol{\gamma}_T, \rho_T) = \Pr(\rho_T | \boldsymbol{\gamma}_T, \rho_{T-1}) \Pr(\boldsymbol{\gamma}_T) \times \Pr(\rho_{T-1} | \boldsymbol{\gamma}_{T-1}, \rho_{T-2}) \Pr(\boldsymbol{\gamma}_{T-1}) \cdots \times \Pr(\rho_2 | \boldsymbol{\gamma}_2, \rho_1) \Pr(\boldsymbol{\gamma}_2) \Pr(\rho_1). \quad (4)$$

We describe $\Pr(\rho_t | \boldsymbol{\gamma}_t, \rho_{t-1})$ shortly, but first provide a definition.

*best, very little trace lefts
of w/e not to be considered*

*P-4
 $\alpha=4$ → all units stays w/e
 to have cluster e/e were at time t-1
 $\alpha=0$ → all units eat
 reallocated
 alpha : more "fixity"
 w/e the clusters, less
 movement
 alpha : lots of tendency
 to change clusters*

Definition 1. We say that partitions ρ_{t-1} and ρ_t are compatible with respect to γ_t , if ρ_t may be obtained from ρ_{t-1} by reallocating items as indicated by γ_t , that is, those items i such that $\gamma_{it} = 0$ for $i = 1, \dots, m$. Note that the compatibility relation is an equivalence relation.

There is a simple way to check if ρ_{t-1} is compatible with ρ_t with respect to γ_t . Let $\mathfrak{R}_t = \{i : \gamma_{it} = 1\}$ be the collection of units that remain fixed when moving from time $t - 1$ to time t , and $\mathfrak{R}_t^C = \{i : \gamma_{it} = 0\}$ is the collection of units that do not. Next denote with $\rho_t^{\mathfrak{R}_t}$ the “reduced” partition at time t that remains after removing all items in \mathfrak{R}_t^C from the subsets of ρ_t . Similarly, let $\rho_{t-1}^{\mathfrak{R}_t}$ be the reduced partition at time $t - 1$ based on γ_t . Then ρ_{t-1} and ρ_t are compatible with respect to γ_t if and only if $\rho_{t-1}^{\mathfrak{R}_t} = \rho_t^{\mathfrak{R}_t}$.

Now, to further characterize $\Pr(\rho_t | \gamma_t, \rho_{t-1})$, let P denote the set of all partitions of m units and let $P_{C_t} = \{\rho_t \in P : \rho_{t-1}^{\mathfrak{R}_t} = \rho_t^{\mathfrak{R}_t}\}$ be the collection of partitions at time t that are compatible with ρ_{t-1} based on γ_t . Then, by construction, $\Pr(\rho_t | \gamma_t, \rho_{t-1})$ is a random partition distribution whose support is P_{C_t} so that

$$\Pr(\rho_t = \lambda | \gamma_t, \rho_{t-1}) = \frac{\Pr(\rho_t = \lambda) I[\lambda \in P_{C_t}]}{\sum_{\lambda'} \Pr(\rho_t = \lambda') I[\lambda' \in P_{C_t}]},$$

where $\Pr(\rho_t = \lambda)$ is the EPPF at the first time point evaluated at λ . Here, and in what follows, $I[A]$ denotes an indicator function with $I[A] = 1$ if statement A is true, and 0 otherwise.

It would be appealing if marginally each of the ρ_t follow the assumed EPPF for ρ_1 , so that the joint probability model for partitions would be stationary. The following proposition establishes this result, which is a consequence of the fact that conditioning on γ_t provides a “reduced” EPPF.

Proposition 1. Let $\rho_1 \sim$ EPPF and $\gamma_1 = \mathbf{0}$. If a joint model for ρ_1, \dots, ρ_T is constructed as described above by introducing γ_t for $t = 2, \dots, T$, then we have that marginally ρ_1, \dots, ρ_T are identically distributed with law coming from the EPPF used to model ρ_1 . Specifically, letting $\rho_{-t} = (\rho_1, \dots, \rho_{t-1}, \rho_{t+1}, \dots, \rho_T)$ and $\gamma = (\gamma_1, \dots, \gamma_T)$, we have that for all $\lambda \in P$,

$$\begin{aligned} \Pr(\rho_t = \lambda) &= \sum_{\rho_{-t} \in P^{\otimes m}} \sum_{\gamma \in \Gamma^{\otimes m}} \Pr(\gamma_1, \rho_1, \dots, \rho_t = \lambda, \dots, \gamma_T, \rho_T) \\ &= \Pr(\rho_1 = \lambda), \end{aligned}$$

where $P^{\otimes m} = P \times P \times \dots \times P$, P a collection of all partitions of m units, $\Gamma^{\otimes m} = \Gamma \times \Gamma \times \dots \times \Gamma$, and Γ a collection of all possible binary vectors of size m .

Proof. See supplementary material. \square

In what follows, we will use $tRPM(\alpha, M)$ to denote our temporal random partition model (4) parameterized by $\alpha_1, \dots, \alpha_T$ and the EPPF in (2). We briefly mention that introducing γ_t is similar in spirit to the approach taken by Caron, Davy, and Doucet (2007); Caron et al. (2017). However, they use γ_t to identify a partial partition at time t that informs how *all* the observational units will be reallocated at time $t + 1$. While this difference may seem inconsequential at first glance, it has drastic ramifications on the type of dependence that exists among the

actual sequence of partitions. This is illustrated in Figures 1 and 2 provided in the Introduction. The sequence of partitions used to create Figure 2 were generated using the $tRPM(\alpha, M)$ with $M = 0.5$. As mentioned in the Introduction, when the main emphasis is on modeling a “smoothly” evolving sequence of random partitions, the temporal dependence displayed in Figure 2 is much more natural than that found in Figure 1.

2.2. Dependence in Partitions

We now further explore how our method models dependence across partitions. To do this, we analyze closeness between partitions ρ_1 and ρ_2 by way of co-clustering of cluster labels (c_{11}, \dots, c_{m1}) and (c_{12}, \dots, c_{m2}) , respectively. We base our exploration on the Rand index which is defined as

$$R(\rho_1, \rho_2) = \frac{a + b}{\binom{m}{2}},$$

where a is the number of pairs (i, j) with $i, j \in [m] = \{1, \dots, m\}$ that simultaneously co-cluster in ρ_1 and ρ_2 and b is the number of such pairs that simultaneously do not co-cluster. Writing $\varphi_{ij} = P(c_{i1} = c_{j1}, c_{i2} = c_{j2}) + P(c_{i1} \neq c_{j1}, c_{i2} \neq c_{j2})$, we note that

$$E[R(\rho_1, \rho_2)] = \binom{m}{2}^{-1} \sum_{1 \leq i < j \leq m} \varphi_{ij}.$$

To provide context to the co-clustering probabilities of cluster labels based on $tRPM(\alpha, M)$, we also consider the model proposed in Caron, Davy, and Doucet (2007). In their approach and assuming $\rho_1 \sim CRP(M)$, each $i \in [m]$ is randomly removed from the partition with probability $1 - \alpha$, and ρ_2 is formed by running an extra $CRP(M)$ process, but starting from an urn that has weights given by the normalized cluster sizes left from the removal process. See details in Caron, Davy, and Doucet (2007). We denote partitions that follow the model of Caron, Davy, and Doucet (2007) as $\rho_1, \rho_2 \sim CAR(\alpha, M)$. For both our method and $CAR(\alpha, M)$, the case that $\alpha = 0$ leads to $\rho_1, \rho_2 \stackrel{iid}{\sim} CRP(M)$, while the largest degree of dependence between ρ_1 and ρ_2 is achieved when $\alpha = 1$. The following proposition characterizes the co-clustering probabilities under our method and $CAR(\alpha, M)$.

Proposition 2. Let $m = T = 2$, so that $E[R(\rho_1, \rho_2)] = \varphi_{12}$.

(a) If $\rho_1, \rho_2 \sim tRPM(\alpha, M)$, where to simplify notation we write $\alpha \equiv \alpha_2$, then

$$\varphi_{12} = \alpha^2 + \frac{(1 + M^2)}{(1 + M)^2} (1 - \alpha)^2.$$

(b) If $\rho_1, \rho_2 \sim CAR(\alpha, M)$ then

$$\varphi_{12} = \left[\frac{6 + 3M + 4M^2 + M^3}{(M + 1)(M + 2)(M + 3)} \right] \alpha^2 + \frac{(1 + M^2)}{(1 + M)^2} (1 - \alpha^2).$$

Proof. See supplementary material. \square

An interesting consequence of Proposition 2 is that we can compute the expected value of the Rand index in the case $\rho_1, \rho_2 \stackrel{iid}{\sim} CRP(M)$.

Corollary 1. If $\rho_1, \rho_2 \stackrel{\text{iid}}{\sim} \text{CRP}(M)$ then for any $m \geq 2$,

$$E[R(\rho_1, \rho_2)] = \frac{(1+M^2)}{(1+M)^2}.$$

Proof. The result follows immediately by noting that the i.i.d. case coincides with tRPM(0, M) and that by exchangeability and independence, $\varphi_{ij} = \varphi_{12}$ for all $1 \leq i < j \leq m$. \square

The result from [Proposition 2](#) (a) shows that, under the tRPM(α, M) model, $\lim_{\alpha \rightarrow 0^+} \varphi_{12} = \frac{(1+M^2)}{(1+M)^2}$, that is, it agrees with $E[R(\rho_1, \rho_2)]$ under the iid case. The same holds as $\alpha \rightarrow 0^+$ under the CAR(α, M) model. Furthermore, for the tRPM(α, M), we get the appealing result that $\lim_{\alpha \rightarrow 1^+} \varphi_{12} = 1$, but the same limit under the CAR(α, M) is a number strictly less than 1 for any $M > 0$. This reveals that the closeness between partitions under the proposed tRPM(α, M), as measured by the φ_{12} quantity, can attain its maximum value of 1, which simply corresponds to the case where none of the units is relocated. The same cannot hold for the CAR(α, M) model because partitions are linked through a latent mechanism rather than directly as in the proposed model. Finally, we conjecture that similar results can be obtained for $m > 2$ but calculations become more involved. The result from [Corollary 1](#) is nevertheless valid for any $m \geq 2$.

2.3. Toy Example to Illustrate Conditional Model

To build intuition regarding the transition from ρ_{t-1} to ρ_t , consider the conditional probabilities in [Equation \(4\)](#) and the very simple scenario of $m = 3$ and $T = 2$. We have that

$$\Pr(\rho_2 | \rho_1) = \sum_{\gamma_2 \in \Gamma} \Pr(\rho_2 | \gamma_2, \rho_1) \Pr(\gamma_2),$$

where again, Γ is the collection of all possible binary 3-tuples and operate under $\rho_1 \sim \text{CRP}(M)$. The conditional probabilities are provided in [Table 1](#), where we set $M = 1$ for simplicity. From [Table 1](#), notice that $\Pr(\rho_2 | \rho_1)$ is a reweighted CRP and that, as $\alpha \rightarrow 0$, partition probabilities correspond to those from the original CRP and, as $\alpha \rightarrow 1$, $\Pr(\rho_2 = \rho_1) \rightarrow 1$. Further, notice that partitions associated with ρ_2 that are more similar to ρ_1 are given larger weight relative to a CRP. For example, given $\rho_1 = \{\{1, 2, 3\}\}$ then $\rho_2 = \{\{1, 2\}, \{3\}\}$ has higher probability than $\rho_2 = \{\{1\}, \{2\}, \{3\}\}$ for any $\alpha > 0$ but have equal probability in a CRP. From this toy example we see that the conditional co-clustering probabilities display dependencies in line with the desire to have partitions evolve gently over time.

2.4. Hierarchical Data Model

Once a partition model is specified, there is tremendous flexibility regarding how to model time (global or cluster-specific) at different levels of a hierarchical model (at the data level, process level, or both). Since we are interested to see how including time in the partition model impacts clustering and model fits, in the simulations of [Section 3](#), we consider a hierarchical model where time only appears in the partition model. In particular, using cluster label notation, we will employ the following hierarchical model

$$\begin{aligned} Y_{it} | \mu_t^*, \sigma_t^{2*}, c_t &\stackrel{\text{ind}}{\sim} N(\mu_{c_{it}t}^*, \sigma_{c_{it}t}^{2*}), i = 1, \dots, m \text{ and } t = 1, \dots, T, \\ (\mu_{jt}^*, \sigma_{jt}^*) | \theta_t, \tau_t^2 &\stackrel{\text{ind}}{\sim} N(\theta_t, \tau_t^2) \times \text{UN}(0, A_\sigma), j = 1, \dots, k_t, \\ (\theta_t, \tau_t) &\stackrel{\text{iid}}{\sim} N(\phi_0, \lambda^2) \times \text{UN}(0, A_\tau), t = 1, \dots, T, \\ (\phi_0, \lambda) &\sim N(m_0, s_0^2) \times \text{UN}(0, A_\lambda), \quad \text{✓-6} \\ \{c_t, \dots, c_T\} &\sim \text{tRPM}(\alpha, M), \text{ with } \alpha_t \stackrel{\text{iid}}{\sim} \text{Beta}(a_\alpha, b_\alpha), \end{aligned} \quad (5)$$

where Y_{it} denotes the response measured on the i th unit at time t , UN denotes a uniform distribution and $A_\sigma, A_\tau, A_\lambda, m_0, s_0^2, a_\alpha, b_\alpha, M$ are user-supplied hyperparameters. The remaining assumptions (e.g., independence across clusters and exchangeability within each cluster) are commonly employed.

2.5. Computation

As the posterior distribution implied by model (5) is not of a known form, we build an algorithm to sample from it. The construction of $\Pr(\rho_1, \dots, \rho_T)$ naturally leads one to consider a Gibbs sampler. In the Gibbs sampler, γ_t will need to be updated in addition to ρ_t (by way of c_t). But the Markovian assumption reduces some of the cost as we only need to consider ρ_{t-1} and ρ_{t+1} when updating ρ_t . Even though each update of ρ_t and γ_t for $t = 1, \dots, T$ needs to be checked for compatibility, it is fairly straightforward to adapt standard algorithms, for example, Algorithm 8 of [Neal \(2000\)](#), with care to make sure that only experimental units with $\gamma_{it} = 0$ are considered when updating c_{it} . Here we provide a general sketch for updating c_{it} and γ_{it} within an MCMC algorithm, with much more detail provided in [Section B](#) of the online supplementary material.

The MCMC algorithm we employ depends on deriving the complete conditionals for ρ_t and γ_t . A key result needed to derive them is provided in the following proposition.

Proposition 3. Based on the construction of a joint probability model as described in [Section 2.1](#), we have

$$\Pr(\rho_t | \gamma_t, \rho_{t-1}) = \begin{cases} \Pr(\rho_t) / \Pr(\rho_t^{\mathfrak{R}_t}) & \text{if } \rho_{t-1}^{\mathfrak{R}_t} = \rho_t^{\mathfrak{R}_t} \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

Table 1. Partition probabilities from the conditional distribution $\Pr(\rho_2 | \rho_1)$ using a CRP EPPF.

(c_1, c_2, c_3)	$\Pr(\rho_1)$	$\Pr(\rho_2 \rho_1 = a)$	$\Pr(\rho_2 \rho_1 = b)$	$\Pr(\rho_2 \rho_1 = c)$	$\Pr(\rho_2 \rho_1 = d)$	$\Pr(\rho_2 \rho_1 = e)$
$a = (1, 1, 1)$	$\frac{2}{6}$	$\frac{2}{6}[1 + 3\alpha^2 - \alpha^3]$	$\frac{2}{6}[1 - \alpha^2]$	$\frac{2}{6}[1 - \alpha^2]$	$\frac{2}{6}[1 - \alpha^2]$	$\frac{2}{6}[1 - 3\alpha^2 + 2\alpha^3]$
$b = (1, 1, 2)$	$\frac{1}{6}$	$\frac{1}{6}[1 - \alpha^2]$	$\frac{1}{6}[1 + 3\alpha^2 + 2\alpha^3]$	$\frac{1}{6}[1 - \alpha^2]$	$\frac{1}{6}[1 - \alpha^2]$	$\frac{1}{6}[1 + \alpha^2 - 2\alpha^3]$
$c = (1, 2, 1)$	$\frac{1}{6}$	$\frac{1}{6}[1 - \alpha^2]$	$\frac{1}{6}[1 - \alpha^2]$	$\frac{1}{6}[1 + 3\alpha^2 + 2\alpha^3]$	$\frac{1}{6}[1 - \alpha^2]$	$\frac{1}{6}[1 + \alpha^2 - 2\alpha^3]$
$d = (1, 2, 2)$	$\frac{1}{6}$	$\frac{1}{6}[1 - \alpha^2]$	$\frac{1}{6}[1 - \alpha^2]$	$\frac{1}{6}[1 - \alpha^2]$	$\frac{1}{6}[1 + 3\alpha^2 + 2\alpha^3]$	$\frac{1}{6}[1 + \alpha^2 - 2\alpha^3]$
$e = (1, 2, 3)$	$\frac{1}{6}$	$\frac{1}{6}[1 - 3\alpha^2 + 2\alpha^3]$	$\frac{1}{6}[1 + \alpha^2 - 2\alpha^3]$	$\frac{1}{6}[1 + \alpha^2 - 2\alpha^3]$	$\frac{1}{6}[1 + \alpha^2 - 2\alpha^3]$	$\frac{1}{6}[1 + 3\alpha^2 + 2\alpha^3]$

Proof. See the supplementary material. \square

When updating γ_{it} in a Gibbs sampler, one can think of removing γ_{it} from $\boldsymbol{\gamma}_t$, and then reinsert it as either a 0 or 1. To this end let $\mathfrak{R}_t^{(-i)} = \mathfrak{R}_t \setminus \{i\}$ and $\mathfrak{R}_t^{(+i)} = \mathfrak{R}_t^{(-i)} \cup \{i\}$ and let $\boldsymbol{\gamma}_{t,+i}$ denote the $\boldsymbol{\gamma}_t$ vector with the i th entry set to 1. Then the full conditional for $\gamma_{it} = 1$, denoted by $\Pr(\gamma_{it} = 1| -)$, is

$$\begin{aligned}\Pr(\gamma_{it} = 1| -) &\propto \Pr(\rho_t | \boldsymbol{\gamma}_{t,+i}, \rho_{t-1}) \Pr(\boldsymbol{\gamma}_{t,+i}) \mathbb{I}[\rho_{t-1}^{\mathfrak{R}_t^{(+i)}} = \rho_t^{\mathfrak{R}_t^{(+i)}}], \\ &\propto \frac{\Pr(\rho_t)}{\Pr(\rho_t^{\mathfrak{R}_t^{(+i)}})} \alpha_t^{\gamma_{it}} \mathbb{I}[\rho_{t-1}^{\mathfrak{R}_t^{(+i)}} = \rho_t^{\mathfrak{R}_t^{(+i)}}],\end{aligned}$$

which results in

$$\begin{aligned}\Pr(\gamma_{it} = 1| -) &= \frac{\alpha_t \Pr(\rho_t^{\mathfrak{R}_t^{(-i)}})}{\alpha_t \Pr(\rho_t^{\mathfrak{R}_t^{(-i)}}) + (1 - \alpha_t) \Pr(\rho_t^{\mathfrak{R}_t^{(+i)}})} \mathbb{I}[\rho_{t-1}^{\mathfrak{R}_t^{(+i)}} = \rho_t^{\mathfrak{R}_t^{(+i)}}].\end{aligned}\quad (7)$$

For a given EPPF that has a closed form (e.g., CRP), it is straightforward to compute $\Pr(\rho_t^{\mathfrak{R}_t^{(-i)}})$ and $\Pr(\rho_t^{\mathfrak{R}_t^{(+i)}})$. If, however, the

$$\Pr(c_{it} = h| -) \propto \begin{cases} N(Y_{it} | \mu_{c_{it}=h,t}^*, \sigma_{c_{it}=h,t}^{2*}) \Pr(c_{it} = h) \mathbb{I}[\rho_{t:c_{it}=h}^{\mathfrak{R}_{t+1}} = \rho_{t+1}^{\mathfrak{R}_{t+1}}] & \text{for } h = 1, \dots, k_t^{-i}, \\ N(Y_{it} | \mu_{new_{h,t}}^*, \sigma_{new_{h,t}}^{2*}) \Pr(c_{it} = h) \mathbb{I}[\rho_{t:c_{it}=h}^{\mathfrak{R}_{t+1}} = \rho_{t+1}^{\mathfrak{R}_{t+1}}] & \text{for } h = k_t^{-i} + 1, \end{cases}\quad (9)$$

where $\mu_{new_{h,t}}^*$ along with $\sigma_{new_{h,t}}^{2*}$ are auxiliary parameters drawn from the prior as in Neal (2000)'s Algorithm 8 (with one auxiliary parameter) and k_t^{-i} is the number of clusters at time t when the i th unit has been removed. Further $N(\cdot | m, s^2)$ denotes a normal density with mean m and variance s^2 . Given ρ_t and $\boldsymbol{\gamma}_t$, the full conditionals of the remaining parameters in model (5) follow standard techniques. A sample can be drawn from the posterior distribution implied by model (5) by iterating through the complete conditionals for $\boldsymbol{\gamma}_t$ and ρ_t and those of other model parameters. See Section B.2 of the online supplementary material for more detail.

In our experience the MCMC algorithm described here and in Section B of the online supplementary material generally behaves well with regards to mixing and convergence. However, applications where $\alpha \approx 1$ can negatively affect the performance of the algorithm. Having a parameter close to the boundary of its support commonly produces computational issues. It is possible to mitigate this by selecting a prior for α that keeps it from its boundary. Alternatively, a specialized algorithm will be needed to accommodate the boundary effect.

3. Simulation Studies

In this section we detail three simulation studies that illustrate different aspects of our modeling approach. In Section C of the online supplemental material, we provide additional simulation results and details regarding a fourth simulation study that considers the performance of our method when the response exhibits spatio-temporal dependence.

EPPF does not have a closed form, then note that (7) can be re-expressed as

$$\Pr(\gamma_{it} = 1| -) = \frac{\alpha_t}{\alpha_t + (1 - \alpha_t) \Pr(\rho_t^{\mathfrak{R}_t^{(+i)}}) / \Pr(\rho_t^{\mathfrak{R}_t^{(-i)}})} \mathbb{I}[\rho_{t-1}^{\mathfrak{R}_t^{(+i)}} = \rho_t^{\mathfrak{R}_t^{(+i)}}].\quad (8)$$

The quantity $\Pr(\rho_t^{\mathfrak{R}_t^{(+i)}}) / \Pr(\rho_t^{\mathfrak{R}_t^{(-i)}})$ is a commonly encountered expression in MCMC methods that employ Neal's Algorithm 8 (Neal 2000). Those same methods can be employed to calculate the desired probabilities. See Section B.1 of the online supplementary material for more detail.

When updating c_{it} note that, within the MCMC algorithm, only those c_{it} for which $\gamma_{it} = 0$ are updated. Thus $\rho_{t-1}^{\mathfrak{R}_t} = \rho_t^{\mathfrak{R}_t}$ by construction. As a result, only compatibility between ρ_t and ρ_{t+1} (i.e., $\rho_t^{\mathfrak{R}_{t+1}} = \rho_{t+1}^{\mathfrak{R}_{t+1}}$) needs to be checked when updating c_{it} . Now letting $\Pr(c_{it} = h) = \Pr(c_{1t}, \dots, c_{it} = h, \dots, c_{mt})$ and denoting the partition based on $\{c_{1t}, \dots, c_{it} = h, \dots, c_{mt}\}$ as $\rho_{t:c_{it}=h} = \{S_{1t}^{-i}, \dots, S_{ht}^{-i} \cup \{i\}, \dots, S_{k_t^{-i}t}^{-i}\}$ where S_{jt}^{-i} denotes the j th cluster at time t with the i th unit removed (note it is possible that $S_{jt}^{-i} = S_{jt}$), the full conditional multinomial probability for c_{it} is provided in (9) where

MCM once
previous file
seem out
model car
now in the
year, / good
one

3.1. Simulation 1: Temporal Dependence in Estimated Partitions

The purpose of the first simulation is to study the accuracy of partition estimates (i.e., $\hat{\rho}_t$) and how much they change over time. (For a discussion of how $\hat{\rho}_t$ is obtained, see below.) In addition, we explore accuracy in estimating $\mu_{it} = \mu_{c_{it}}$ and α_t . To this end, we considered model (5) as a data-generating mechanism to create 100 datasets with 50 observations at five time points. We used tRPM(α, M) with $\alpha_t = \alpha$ for all t and $M = 1$. We generate synthetic datasets under $\alpha \in \{0, 0.1, 0.25, 0.5, 0.75, 0.9, 0.999\}$. For all i and t , we set $\sigma_{c_{it}}^{2*} = \sigma^2 = 1$, $\tau^2 = 25$, and $\theta_t = 0$.

To each synthetic data set we fit model (5) using the MCMC algorithm detailed in Section 2.5 by collecting 10,000 iterates and discarding the first 5000 as burn-in and thinning by 5 (resulting in 1000 MCMC samples). As prior parameters we used $A_\sigma = 5$, $A_\tau = 10$, $A_\lambda = 10$, $m_0 = 0$, $S_0^2 = 100$, $a_\alpha = b_\alpha = M = 1$. For simplicity we set $\alpha_t = \alpha$ for all t . All partition point estimates were estimated using the method in the `salso` R package (Dahl, Johnson, and Müller 2020) with the binder loss function (Binder 1978). To measure similarity between partitions, we employed the ARI (Rand 1971; Hubert and Arabie 1985) and we used WAIC (Gelman, Hwang, and Vehtari 2014) to measure model fit.

Table 2 displays the lagged 1 and 4 ARI as a function of α . As expected, for both lags, the ARI increases as α increases. Also as expected, lagged 4 ARI increases less as a function of α compared to the lagged 1 ARI. Note that on average the lagged

greatest
metrows

Table 2. ARI when comparing $\hat{\rho}_1$ to $\hat{\rho}_2$ and $\hat{\rho}_1$ to $\hat{\rho}_5$.

	ARI($\hat{\rho}_1, \hat{\rho}_2$)	ARI($\hat{\rho}_1, \hat{\rho}_5$)	Coverage		WAIC	
			α	μ_{it}	tRPM	CRP
$\alpha = 0.0$	0.05 (0.02)	0.01 (0.01)	0.00 (0.00)	0.93 (0.01)	896	892
$\alpha = 0.1$	0.04 (0.01)	0.00 (0.01)	0.96 (0.02)	0.92 (0.01)	910	907
$\alpha = 0.25$	0.19 (0.03)	0.02 (0.02)	0.90 (0.03)	0.91 (0.01)	890	891
$\alpha = 0.5$	0.44 (0.02)	0.04 (0.01)	0.82 (0.04)	0.91 (0.01)	881	903
$\alpha = 0.75$	0.67 (0.02)	0.23 (0.02)	0.89 (0.03)	0.92 (0.01)	822	893
$\alpha = 0.9$	0.87 (0.01)	0.55 (0.02)	0.90 (0.03)	0.90 (0.01)	816	896
$\alpha = 0.9999$	0.97 (0.01)	0.93 (0.01)	0.58 (0.05)	0.90 (0.01)	795	888

NOTES: ARI(., .) denotes the adjusted Rand index as a function of two partitions. Coverage rates for α and μ_{it} and model fit metrics for tRPM(α, M) and CRP(M). These values are averaged over the 100 generated data sets. The values in parenthesis are Monte Carlo standard errors. Smaller values of WAIC indicate better fit.

1 ARI for $\alpha = 0.1$ is smaller than that for $\alpha = 0$. This is because the variability associated with lagged 1 ARI when $\alpha = 0$ is much larger than when $\alpha > 0$, producing a few lagged ARI values that are large. The median of the lagged ARI values increase as a function of α monotonically. Since ARI isn't necessarily a monotone function of α , care must be taken when interpreting α as a dependence parameter based on ARI.

To study the ability to recover μ_{it} and α , 95% credible intervals for each were computed and coverage was estimated. Results are provided in Table 2. Notice that coverage for α is low when the true α is at or near the boundary (e.g., $\alpha \in \{0, 0.9999\}$) which is to be expected. The coverage associated with μ_{it} is close to the nominal rate regardless of the value of α . Therefore, temporal dependence in the partition model does not adversely impact the ability to estimate individual means.

Lastly, to compare model fit when using tRPM(α, M) as the RPM in model (5) relative to $\rho_t \stackrel{iid}{\sim} \text{CRP}(M)$, we calculated the WAIC for each data set when fitting model (5) under both RPMs. Results are provided in Table 2 where each entry is an average WAIC value over all 100 datasets. Notice that, when the independent partitions were used to generate data (i.e., $\alpha = 0$), modeling partitions independently produces slightly better model fit as would be expected. But even if relatively weak temporal dependence exists among the sequence of partitions, there are gains in modeling the sequence of partitions with tRPM(α, M), with gains becoming substantial as α increases.

The upshot from this simulation study is that lagged partition estimates when employing tRPM(α, M) display intuitive behavior in that similarity between partition estimates decreases as lag increases. In addition, employing the tRPM(α, M) partition model does not negatively impact parameter estimation and produces improved model fits when dependence is present in the sequence of partitions and a minimal cost in model fit when it is not.

3.2. Simulation 2: Induced Correlation at the Response Level

A potential benefit of developing a joint model for partitions is the ability to accommodate temporal dependence that may exist between Y_{it} and Y_{it+1} . To study this, we conducted a small Monte Carlo simulation study that is comprised of sampling repeatedly from the tRPM(α, M) using the computational approach of Section 2.5. Once the partition is generated, the temporal dependence among the Y_i depends on specific model

choices for μ_{jt}^* . Here we use $\mu_{jt}^* \sim N(\phi_1 \mu_{jt-1}^*, \tau^2(1 - \phi_1^2))$ for $t > 2$, $j = 1, \dots, k_t$, and $|\phi_1| \leq 1$. For $t = 1$ we use $\mu_{j1}^* \sim N(0, \tau^2)$ and if $k_{t+1} > k_t$ new μ_{jt+1}^* values are drawn from $N(0, \tau^2)$. Now setting $m = 25$, $T = 10$, $\tau = 10$, and $\sigma = 1$, 100 datasets were generated for $\phi_1 \in \{0, 0.25, 0.5, 0.75, 0.9, 1\}$. For each data set generated, the lagged auto-correlations among Y_i were computed for $i = 1, \dots, m$. The results found in Figure 3 are the lagged auto-correlations averaged over the m units for $\alpha \in \{0, 0.25, 0.5, 0.75, 0.9\}$.

As can be seen in Figure 3, when partitions are independent (i.e., $\alpha = 0$), no correlation propagates to the data level. The same can be said if atoms are iid (i.e., $\phi_1 = 0$). As the temporal dependence among μ_{jt}^* increases (i.e., ϕ_1 increases), there is stronger temporal dependence among Y_{i1}, \dots, Y_{iT} . Notice further that this dependence persists longer in time as α increases as one would expect.

3.3. Simulation 3: AR(1)-type synthetic data

In our final simulation experiment, we consider data generated from an AR(1) process. To create synthetic datasets for the i th unit, we employ the following as a data generating-mechanism

$$Y_{it} = \mu_{ci} + \omega Y_{it-1} + \epsilon_{it}, \text{ for } i = 1, \dots, m, \text{ and } t = 1, \dots, T,$$

where $|\omega| < 1$ and $\epsilon_{it} \sim N(0, \nu^2)$. We consider synthetic datasets with four clusters so that $c_{it} \in \{1, 2, 3, 4\}$ corresponding to $\mu \in \{-2, 0, 2, 4\}$. The four clusters are formed by dividing the $m = 100$ units into equal groups of 25. At time points $t = 2, \dots, T$, four units from each cluster are shifted to other clusters in a systematic way so that clusters change over time. An example of the type of data this procedures creates can be seen in Figure S.1 of the supplementary material. Data are generated using the function `arima.sim` in R (R Core Team 2020) under $\omega \in \{0, 0.1, 0.25, 0.5, 0.75, 0.9\}$, $\nu^2 \in \{0.5^2, 1^2\}$, and $T \in \{5, 10\}$. A total of 100 datasets are created under each scenario (totaling 28) and to each we fit the following competing models.

1. A weighted-dependent Dirichlet process (wddp) described in Quintana et al. (2020) and chapter 4.4.4 of Müller et al. (2015). This model incorporates time in the weights of a DDP. As such, we fit this model to a concatenated version of the data (Y_i, t_i) for $i = 1, \dots, Tm$. Specific details of this procedure are provided in Section C.2 of the supplementary material.
2. A linear dependent Dirichlet process (lddp) described in Quintana et al. (2020) and chapter 4.4.2 of Müller et al.

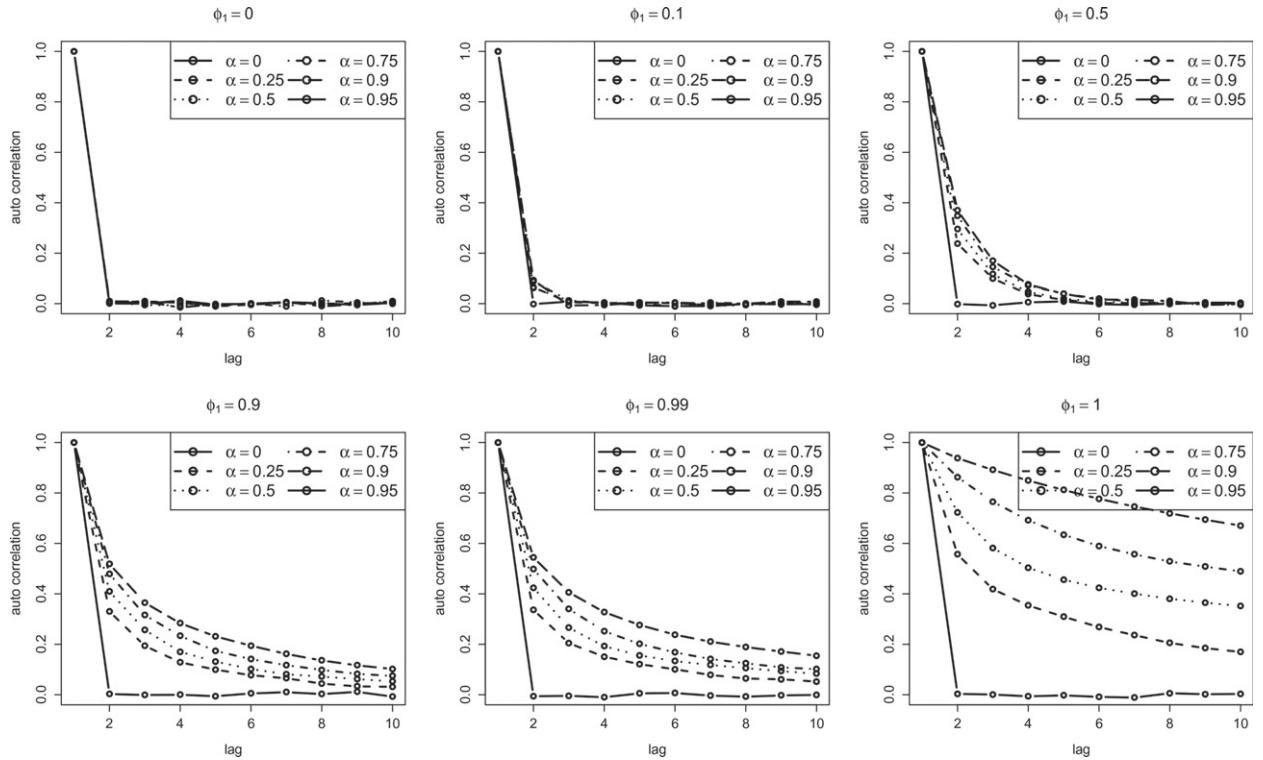


Figure 3. Lagged auto-correlations among the (Y_{i1}, \dots, Y_{iT}) when modeling μ_{jt}^* with an AR(1) type structure.

(2015). This model incorporates time in the atoms of a DDP. As in the wddp, to this model we concatenated version of the data (Y_i, t_i) for $i = 1, \dots, T_m$. Specific details of this procedure are also provided in Section C.2 of the supplementary material.

3. A Griffiths-Milne dependent Dirichlet process (gmddp) mixture. This model was fit using the BNPmix package in R (Corradin, Canale, and Nipoti 2020). See Corradin, Canale, and Nipoti (2021) for specific model details.
4. A temporally independent CRP(M) model (ind_crp). This procedure corresponds to $\alpha = 0$ and is a special case of our approach and that proposed in Caron, Davy, and Doucet (2007). The exact details of this model are also provided in Section C.2 of the supplementary material.
5. A temporally static CRP(M) model (static_crp). This procedure corresponds to $\alpha = 1$ and is also a special case of the model detailed in Caron, Davy, and Doucet (2007). Like the wddp and lddp, this procedure is fit to a concatenated version of the data (Y_i, t_i) for $i = 1, \dots, T_m$. See Section C.2 of the supplementary material for more details.

Results from this simulation study are presented in Figure 4. The left plot in the figure displays the WAIC model fit metric for each procedure averaged over the 100 synthetic datasets, while the right plot displays the ARI value. The ARI values were produced by calculating ARI for each MCMC sample from the posterior distribution of ρ_t . This was done separately for each $t = 1, \dots, T$ and then averaged across time and MCMC samples. From the left plot, our method (drpm) is superior to all other methods in terms of the model fit metric WAIC, save for the lddp method. Against the lddp method, the drpm produces better fits when the data are generated with high auto-

correlation and larger data noise. This is to be expected as the drpm only incorporates temporal information in the prior on partitions while the lddp includes this information in the atoms (i.e., likelihood). That said, in terms of partition estimation, the drpm easily outperforms all other competitors in terms of ARI. The fact that drpm outperformed wddp and lddp in terms of ARI is not surprising as the later methods incorporate time differently compared to the drpm. However, the drpm and gmddp methods both treat time similarly, but the former from a partition perspective while the latter from a random probability measure perspective. This highlights the fact that when partitions are of interest, modeling them directly provides benefit.

4. Application

In this section we apply our method to a real-world data set coming from the field of environmental science. A second application in education is provided in Section D.1 of the supplementary material. As mentioned previously, once a partition model is specified, there is quite a bit of flexibility regarding how (or if) temporal dependence is incorporated in other parts of a hierarchical model. To illustrate this, we incorporate temporal dependence in three places of the hierarchical model we construct.

As part of preliminary exploratory data analysis (not shown), we examined serial dependence for each experimental unit (monitoring station), and concluded that they all exhibited a particular type of temporal dependence. Because of this, we introduce a unit-specific temporal dependence parameter $|\eta_{1i}| \leq 1$ and model observations from a single unit over time (Y_{i1}, \dots, Y_{iT}) with an AR(1) structure. In addition, motivated

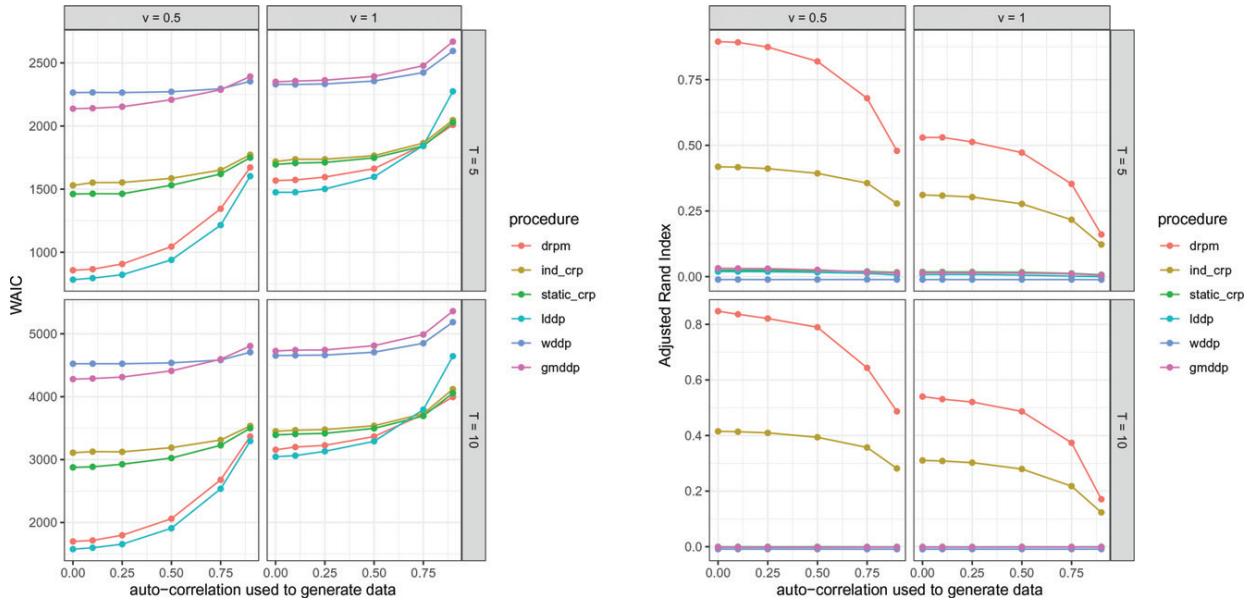


Figure 4. Results from simulation study with data that contains AR(1)-type temporal correlation. The left plot corresponds to results for the model fit metric WAIC and the right plot displays results associated with ARI.

by a desire for parsimony, we employed a Laplace prior for η_{1i} . Finally, to permit the temporal dependence in the partition model to propagate through the hierarchical model, we assume an AR(1) structure for the θ_t 's. The full hierarchical model is detailed in Equation (10).

$$\begin{aligned}
 Y_{it} | Y_{it-1}, \mu_t^*, \sigma_t^{2*}, \eta, c_t &\sim \text{ind } N(\mu_{c_{it}}^* + \eta_{1i} Y_{it-1}, \sigma_{c_{it}}^{2*}(1 - \eta_{1i}^2)), \\
 Y_{i1} &\sim \text{ind } N(\mu_{c_{i1}}^*, \sigma_{c_{i1}}^{2*}), \\
 \xi_i = \text{Logit}(0.5(\eta_{1i} + 1)) &\sim \text{iid Laplace}(a, b), \\
 (\mu_{jt}^*, \sigma_{jt}^*) &\sim \text{ind } N(\theta_t, \tau_t^2) \times \text{UN}(0, A_\sigma), \\
 \theta_t | \theta_{t-1} &\sim N((1 - \phi_1)\phi_0 + \phi_1\theta_{t-1}, \lambda^2(1 - \phi_1^2)), \\
 (\phi_1, \tau_t) &\sim N(\phi_0, \lambda^2) \times \text{UN}(0, A_\tau), \\
 (\phi_0, \phi_1, \lambda) &\sim N(m_0, s_0^2) \times \text{UN}(-1, 1) \times \text{UN}(0, A_\lambda), \\
 \{c_t, \dots, c_T\} &\sim \text{tRPM}(\alpha, M), \text{ with } \alpha_t \stackrel{\text{iid}}{\sim} \text{Beta}(a_\alpha, b_\alpha),
 \end{aligned}
 \tag{10}$$

where all Roman letters correspond to parameters that are user supplied. There are a number of special cases embedded in our hierarchical model. For example, $\eta_{1i} = 0$ for all i results in conditionally independent observations. Further, $\phi_1 = 0$ results in "independent atoms" and $\alpha_t = 0$ for all t in independent partitions over time. The model (5) used in the simulation studies is a special case of (10), obtained by setting $\phi_1 = 0$ and $\eta_{1i} = 0$ for all i .

4.1. Rural Background PM₁₀ Data Application

The rural background PM₁₀ data is taken from the European air quality database. These data are comprised of the daily measurements of particulate matter with a diameter less than 10 μm from rural background stations in Germany and are publicly available in the `gstat` package (Gräler, Pebesma, and Heuvelink 2016) found on CRAN in R (R Core Team 2020). We focus on average monthly PM₁₀ measures from the year 2005 (i.e., $T = 12$). Of the 69 stations, 9 were removed because of missing values.

Table 3. PM₁₀ data: Results of model fitting. The bold font identifies best model fits in terms of LPML and WAIC. Higher values for LPML indicate better fit while lower values for WAIC indicate better fit.

Temporal Dependence In				
Likelihood (η_{1i})	Atoms (ϕ_1)	Partition (α_t)	LPML	WAIC
No	No	No	-1814	3683
No	No	Yes	-1656	3031
No	Yes	No	-1752	3539
No	Yes	Yes	-1644	3271
Yes	No	No	-1704	3554
Yes	No	Yes	-1578	3186
Yes	Yes	No	-1706	3544
Yes	Yes	Yes	-1586	3153

We fit the hierarchical model (10) to these data and consider all the possible special cases (i.e., $\eta_{1i} = 0$ or not, $\phi_1 = 0$ or not, $\alpha_t = 0$ or not). This resulted in 8 total models that were fit by collecting 1000 MCMC iterates after discarding the first 10,000 as burn-in and thinning by 40 (i.e., 50,000 total MCMC samples were collected). Running the algorithm for 50,000 samples on a laptop with 16GB of RAM took between 1 and 2.5 minutes. We use the following prior values: $A_\sigma = 10$, $A_\tau = A_\lambda = 5$, $m_0 = 0$, $s_0^2 = 100$, $a = 0$, $b = 1$, and $a_\alpha = b_\alpha = 2$. The prior for α_t was specified to encourage α from approaching 1. The WAIC and log pseudo-marginal likelihood (LPML) for each model are presented in Table 3. To improve the computational stability of the LPML, for each model fit, the MCMC iterates that correspond to 0.5% of the smallest likelihood values were not included in the calculation of LPML.

First we note that among all the model fits, employing a variant of tRPM(α, M) (i.e., rows with "Yes" in the "Partition" column) improves model fit. The best performing model in terms of WAIC includes temporal dependence in the partition model only, while that for LPML includes temporal dependence in the partition model and in the likelihood.

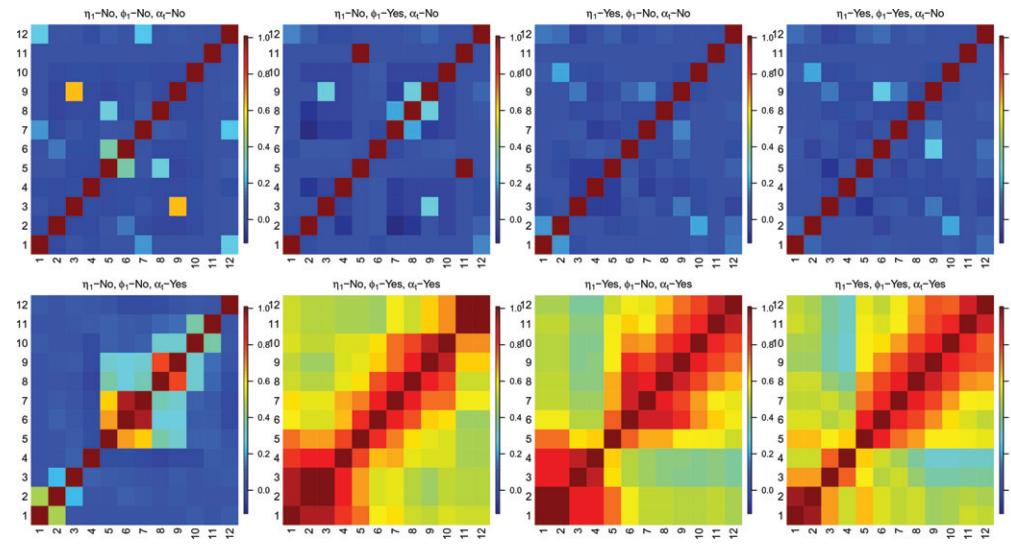


Figure 5. PM₁₀ data. Each figure is a summary of the lagged ARI values corresponding to the 8 models in Table 3. At each time point, the partition was estimated using the `salso` function in the `salso` R package (Dahl, Johnson, and Müller 2020) based on binder loss.

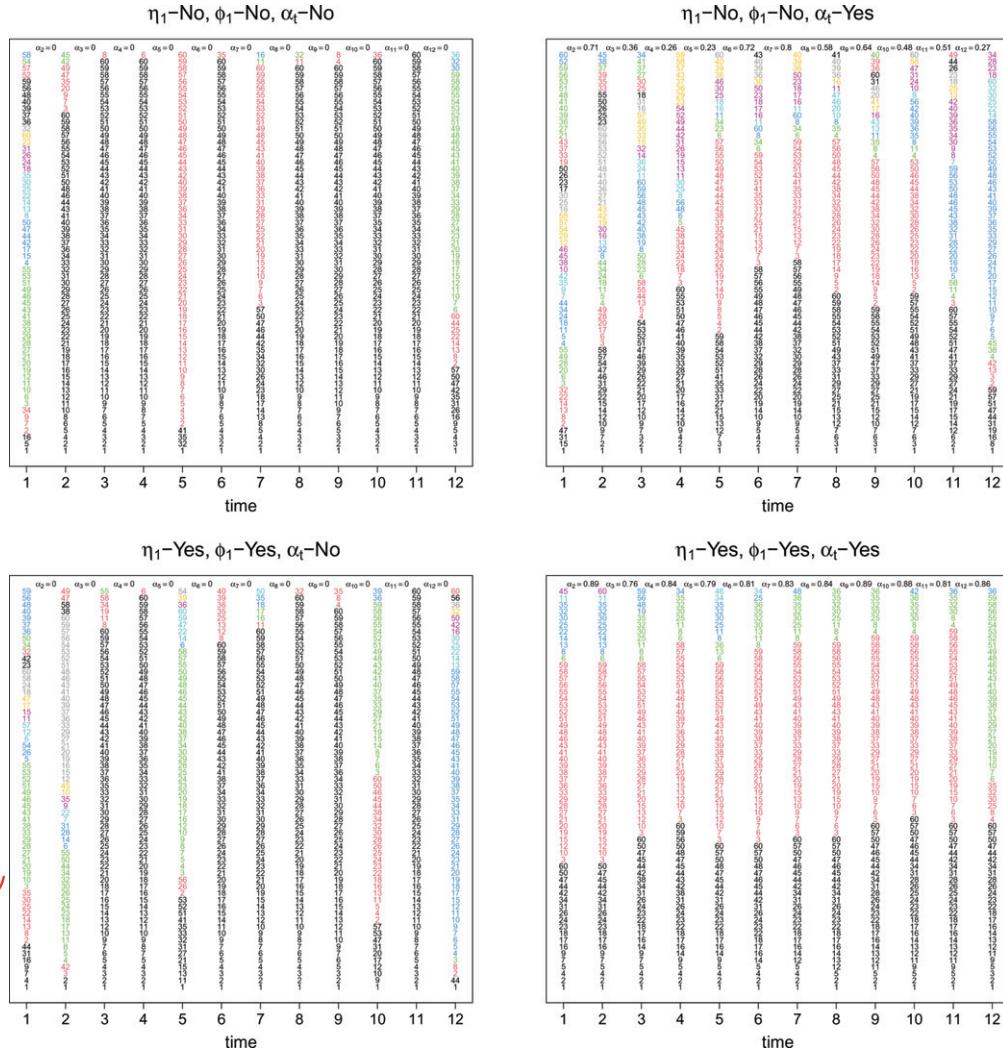


Figure 6. Each plot displays the estimated partition at each time point. Plots in the left column are based on $\rho_t \stackrel{\text{iid}}{\sim} \text{CRP}(M)$ while those in the right column are based on $(\rho_1, \dots, \rho_T) \sim \text{tRPM}(\alpha, M)$. Clusters are highlighted by color and each number corresponds to a specific monitoring station.

Now focusing on partition inference, we provide Figure 5. This figure displays the lagged ARI values for each of the 8 models. Notice that when partitions are modeled independently

(first third row of Figure 5) then partitions evolve over time quite erratically in the sense that the cluster configuration can change dramatically from one time point to the next. However, when

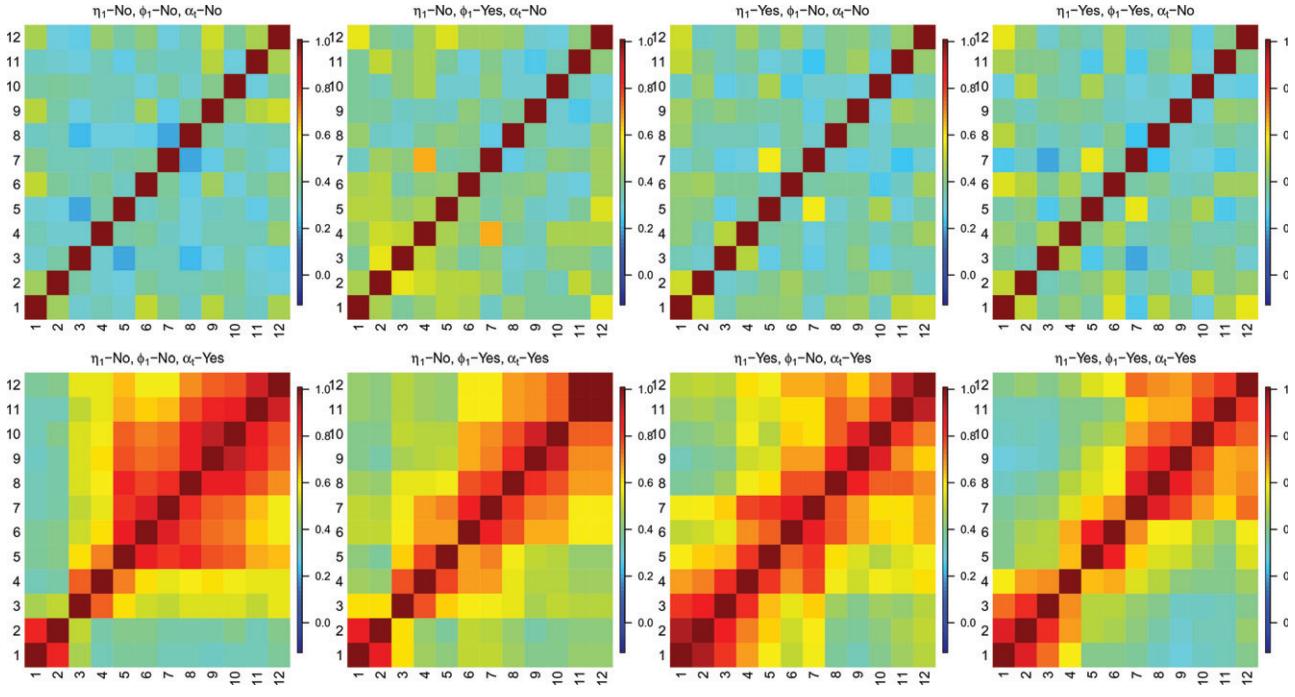


Figure 7. PM₁₀ data. Each figure is a summary of the lagged ARI values corresponding to the 8 models that are detailed in Table 3 except that space is also included in the dependent random partition model. At each time point the partition was estimated using the `salso` function in the `salso` R package (Dahl, Johnson, and Müller 2020) based on binder loss.

employing tRPM(α, M) (second row of Figure 5) the partitions seem to evolve much more “smoothly” as there is less drastic changes in cluster configuration. In fact the model that produces the best model fit metrics (right most plot of second row) seems to produce partitions that change quite gently over time as desired.

Finally, we provide Figure 6 as a means to visualize how estimated partitions based on our joint partition model evolve over time relative to modeling partitions with an iid model. Each plot in Figure 6 displays the estimated partition at each time point. Each color represents a cluster and each number corresponds to a specific measuring station. The plots illustrate the sequential nature of cluster forming with the first cluster always containing the first measuring station, the next cluster is formed by the first station not included in the first cluster and so on. The plots in the right column correspond to using tRPM(α, M) to jointly model partitions while those on the left employ $\rho_t \stackrel{\text{iid}}{\sim} \text{CRP}(M)$. It is evident from Figure 6 that from one time point to the next that partitions based on our construction evolve much more gently over time. This more closely mimics how PM₁₀ measurements would evolve as a function of time compared to the *iid* model.

4.2. Extensions to the Joint Partition Model

Based on our generic joint model construction, it is straightforward to incorporate other information in the partition model such as covariates. For example, in the data application of 4.1 each monitoring station’s spatial coordinates were recorded. Incorporating spatial dependence in our analysis of the PM₁₀ data can be easily accommodated via the EPPF in our construction. This would result in spatially informed clusters that

evolve over time. If the spatially referenced EPPF preserves sample size consistency, then Proposition 1 still holds. One such EPPF is part of the spatial product partition model (sPPM) class developed in Page and Quintana (2016). To illustrate the ease of incorporating space in our model construction, here we model the PM₁₀ data using model (10) but rather than use the tRPM(α, M) to model the sequence of partitions, we use a version of our dependent partition model that employs the sPPM.

In order to introduce the sPPM, let s_i denote the spatial coordinates of the i th item (note that these coordinates do not change over time) and let s_{jt}^* be the subset of spatial coordinates that belong to the j th cluster at time t . Then we express the EPPF of the t th partition with the following product form

$$\Pr(\rho_t | v_0, M) \propto \prod_{j=1}^{k_t} c(S_{jt}|M) g(s_{jt}^* | v_0). \quad (11)$$

Here, $c(\cdot|M) \geq 0$ is called the cohesion and is a set function that produces cluster weights a priori. We consider the cohesion $c(S_{jt}|M) = M \times (|S_{jt}| - 1)!$ as it has connections with the CRP making this version of the sPPM a type of spatially re-weighted CRP. The similarity function $g(\cdot|v_0)$ is a set function parametrized by v_0 that measures the compactness of the spatial coordinates in s_{jt}^* producing higher values if the spatial coordinates in s_{jt}^* are close to each other. Not all similarity functions preserve sample size consistency so to ensure this, after standardizing spatial locations, we employ

$$g(s_{jt}^* | v_0) = \int \prod_{i \in S_{jt}} N(s_i | \mathbf{m}, V) NIW(\mathbf{m}, V | \mathbf{0}, 1, v_0, I) d\mathbf{m} dV, \quad (12)$$

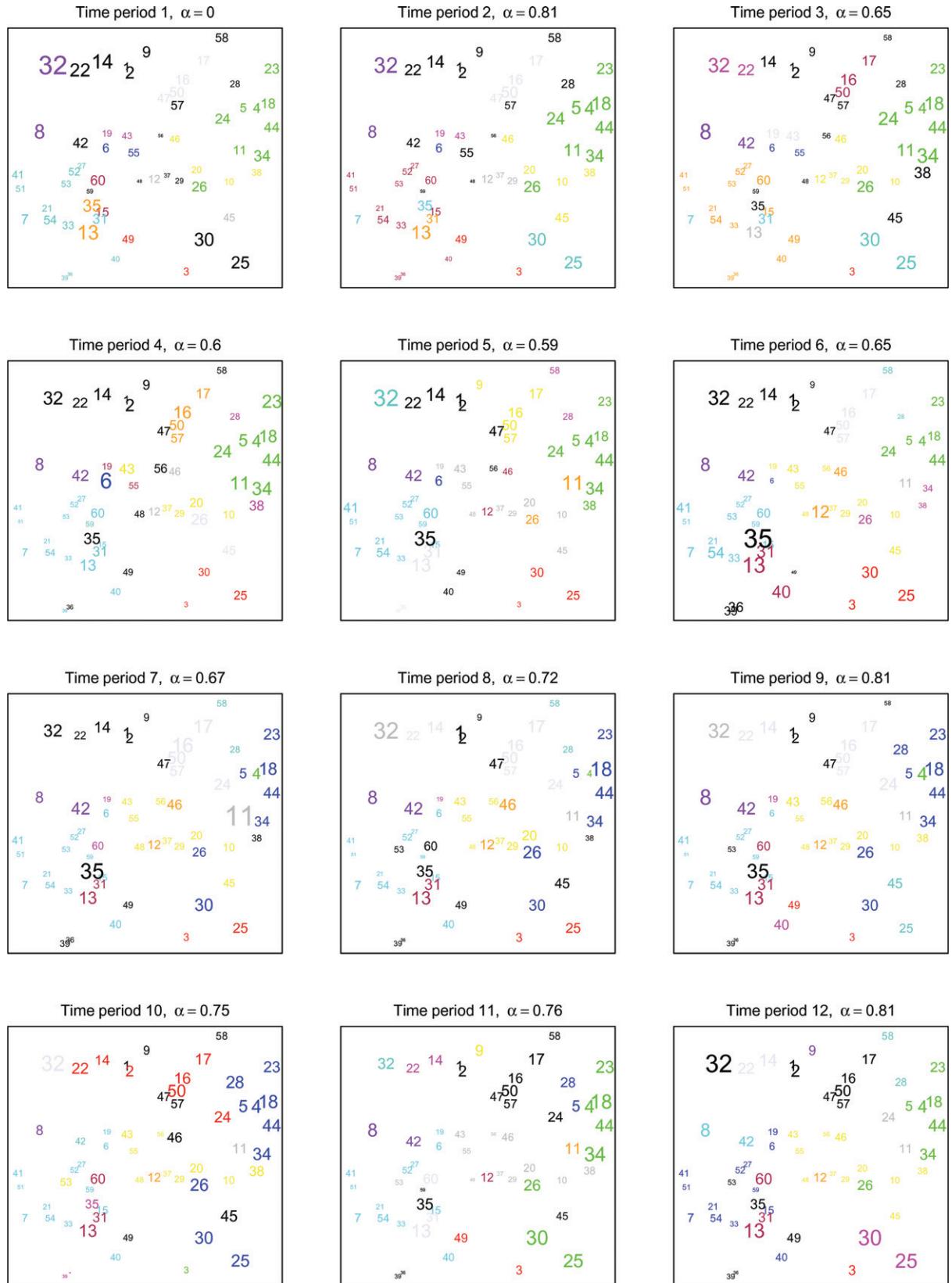


Figure 8. PM₁₀ data. Graphical display for spatially referenced estimated partitions for each time point based on the model that produced best fit (space in the partition model and temporal dependence in all levels of the model). The size of a point is proportional to the PM₁₀ measured at that station. Clusters are identified by color. Each monitoring station is labeled by the same numbers as in Figure 6. At each time point, the partition was estimated using the `salso` function in the `salso` R package (Dahl, Johnson, and Müller 2020) based on binder loss.

where $N(\cdot|\mathbf{m}, \mathbf{V})$ denotes a bivariate normal density and $NIW(\cdot, \cdot|\mathbf{0}, 1, v_0, \mathbf{I})$ a normal-inverse-Wishart density with

mean $\mathbf{0}$, scale equal to 1, inverse scale matrix equal to \mathbf{I} , and v_0 being the user-supplied degrees of freedom. Note that larger

longer

values of ν_0 increase spatial influence on partition probabilities. For more details on why this formulation preserves sample size consistency, see Müller, Quintana, and Rosner (2011) and Quintana, Loschi, and Page (2018). For more information regarding the impact of ν_0 on product form of the partition model, see Page and Quintana (2016, 2018). We will use stRPM(α , ν_0 , M) to denote our spatio-temporal random partition model (4) parameterized by $\alpha_1, \dots, \alpha_T$ and EPPF detailed in Equations (11) and (12).

As in the previous section, we fit model (10) to the PM₁₀ data but replacing tRPM(α , M) with the stRPM(α , ν_0 , M). Also as before, we consider all the possible special cases of the model (i.e., $\eta_{1i} = 0$ or not, $\phi_1 = 0$ or not, $\alpha_t = 0$ or not). This resulted in 8 total models that were fit by collecting 1000 MCMC iterates after discarding the first 10,000 as burn-in and thinning by 40. Fitting each of the eight models based on the stRPM(α , ν_0 , M) took between 20 and 57 minutes. The prior values employed were the same as before with the addition of setting $\nu_0 = 5$ which places in the partition prior moderate weight on spatial locations. *RD-43*

Incorporating space in the partition model seems to provide benefit in terms of model fit as the LPML and WAIC values associated with the model that includes space in the partition model and temporal dependence in all levels of the model fits best in terms of LPML with -1487 compared to -1586 listed in Table 3 and temporal dependence in partition and atoms fits best with regards to WAIC with 3140 compared to 3271 listed in Table 3. Additionally, Figure 7 displays the lagged ARI values for the 8 models that include space. As in Figure 5 when there is no temporal dependence in the partition model, then the estimated partitions exhibit no temporal dependence. However, when time and space are included then there is clear dependence between partitions as a function of time. Further, the dependence between partitions seems to decay faster with space and time included in the partition model compared to just time.

Finally, we provide Figure 8 which displays the estimated spatially referenced partitions at each time point based on the model that achieved the best fit (space in the partition model and temporal dependence in all levels of the model). The size of each point in the figure is proportional to the PM₁₀ measured at the particular station and each color depicts. To make connections with Figure 6 each monitoring station is labeled with the same number as before. Notice that there are clear similarities from one time point to the next for most months. That said, there are two time periods for which changes in the PM₁₀ are more drastic relative to the previous time period (e.g., August to September). In these months the estimated α_t is quite a bit smaller and as a result, the estimated partitions are more different.

5. Conclusions

We developed a joint probability model for a sequence of partitions that explicitly considers temporal dependence among the partitions. Further, we showed that our methodology is capable of accommodating partitions that evolve slowly over time in that the ARI between estimated partitions decays as the lag in time increases. We also showed that if partitions are indeed independent over time, then employing our joint parti-

tion prior regardless results in a minimal cost in terms of model fit.

The predictive nature of the temporal prior on a sequence of random partitions we have presented has a first-order Markovian structure. Various extensions can be considered, such as adding higher order dependence across time or dependence in baseline covariates. All of these cases would build on our constructive definition, as extra refinements of the basic idea of carrying smooth transitions on time (or time and space). Lastly, the Markovian structure could be exploited to carry out predictive inference as well.

Acknowledgments

The first author gratefully acknowledges support from the Basque Government through the BERC 2018-2021 program, by the Spanish Ministry of Science, Innovation and Universities through BCAM Severo Ochoa accreditation SEV-2017-0718. The second author is supported by the grant FONDECYT 1180034 and by ANID - Millennium Science Initiative Program - NCN17_059.

Supplementary Material

Online supplementary material file that contains proofs of propositions, computation details, additional simulation, and application results. *R-package*: The R-package drpm contains the function drpm that fits all models described in the article.

References

- Antoniano-Villalobos, I., and Walker, S. G. (2016), “A Nonparametric Model for Stationary Time Series,” *Journal of Time Series Analysis*, 37, 126–142. [614]
- Binder, D. A. (1978), “Bayesian Cluster Analysis,” *Biometrika*, 65, 31–38. [619]
- Caron, F., Davy, M., and Doucet, A. (2007), “Generalized Polya Urn for Time-Varying Dirichlet Process Mixtures,” in *Proceedings of the Twenty-Third Conference on Uncertainty in Artificial Intelligence (UAI'07)*, Arlington, VA: AUAI Press. Available at: <http://dl.acm.org/citation.cfm?id=3020488.3020493>. [614,617,621]
- Caron, F., Neiswanger, W., Wood, F., Doucet, A., and Davy, M. (2017), “Generalized Pólya Urn for Time-Varying Pitman-Yor Processes,” *Journal of Machine Learning Research*, 18, 1–32. Available at: <http://jmlr.org/papers/v18/10-231.html>. [614,615,617]
- Corradin, R., Canale, A., and Nipoti, B. (2020), “BNPmix: Bayesian Nonparametric Mixture Models.” Available at: <https://CRAN.R-project.org/package=BNPmix>. R package version 0.2.6. [621]
- (2021), “BNPmix: An R Package for Bayesian Nonparametric Modelling Via Pitman–Yor Mixtures,” *Journal of Statistical Software*. [621]
- Dahl, D. B., Johnson, D. J., and Müller, P. (2020), *salso: Search Algorithms and Loss Functions for Bayesian Clustering*. Available at: <https://CRAN.R-project.org/package=salso>. R package version 0.2.5. [619,623,624,625]
- De Blasi, P., Favaro, S., Lijoi, A., Mena, R. H., Prünster, I., and Ruggiero, M. (2015), “Are Gibbs-Type Priors the Most Natural Generalization of the Dirichlet Process?” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37, 212–229. [616]
- De Iorio, M., Favaro, S., Guglielmi, A., and Ye, L. (2019), “Bayesian Nonparametric Temporal Dynamic Clustering Via Autoregressive Dirichlet Priors,” arXiv:1910.10443. [614]
- DeYoreo, M. and Kottas, A. (2018), “Modeling for Dynamic Ordinal Regression Relationships: An Application to Estimating Maturity of Rockfish in California,” *Journal of the American Statistical Association*, 113, 68–80. Available at: <https://doi.org/10.1080/01621459.2017.1328357>. [614]
- Gelman, A., Hwang, J., and Vehtari, A. (2014), “Understanding Predictive Information Criteria for Bayesian Models,” *Statistics and Computing*, 24, 997–1016. [619]

General idea: we have
n distinct locations s_1, \dots, s_m
($s_i = (\text{lat}, \text{long})$ etc)

the goal is to
define a model for the distribution
of s_i into $km(t)$ bins

so we define $\rho_m = \{s_{i1}, \dots, s_{im}\}$? the clusters
definition, i.e. we say \Leftrightarrow location s_i is
in cluster a .
In symbolic form, let c_1, \dots, c_m write
 $c_a = \{s_i\}$ w.r.t. a .

then the idea can be as

$$P_{\rho_m}(\rho) = P(\rho = \rho_m) \propto \prod_{i=1}^{km} C(s_i, c_a)$$

we call this
model / law as
sPPM (symbolic
product parti-
tion model)

this is a
clustering operation:
names low level
these selected
values of s_i would be
clustered together
as near

no model on m w.r.t.
will become $C(s_a, t)$?

then model the process
clustering, depend
location s_a

the set of locations
that $s_i \in s_a$
like if $s_a = s_{i1}, s_{i2}, \dots, s_{im}$
 $c_a = \{s_{i1}, s_{i2}, \dots, s_{im}\}$

Old Part

then the idea is to give each s_i in ρ_m a
low / bivariate value w.r.t. cluster dependent

$$s_i = r(c_i) | \underline{s}_i, c_i \stackrel{\perp}{\sim} f(\underline{s}_{ci}) \quad \text{for } w = u \rightarrow m$$

$\underline{s}_{ci} \approx \underline{s}_i$ and $f(\cdot) = g_0$ for $l = e \rightarrow km$

no actual error
are due to $\rho_m(t)$ as $km = |\rho_m(t)|$

$\{\underline{s}_{ci}\}_{i=1}^m \sim sPPM$

page 3
rec 46

and we can include consume with

let $s_{w,t}$ the values shared (symbolic sense)
let $s_{w,t}$ the consume w.r.t. time
let β^1, \dots, β^m the cluster specific regression terms

the \underline{s}_{ci} of
below

(most, low car we have,
a normal value be near,
i.e., low will be lower influence
high later, if the # km varies
with time, if the β_j are time as not
the β_j are time "will be
the more" will be

then we model

$$s_i = r(c_i) | s_{w,t}, c_i(t), \beta^1, \dots, \beta^m \stackrel{\perp}{\sim} N(\underline{s}_{ci}^\top \beta^1, \dots, \beta^m, \sigma^2)$$

in all this temporal extension
test model, the more

there there are more more examples of models.
note \rightarrow for the ρ_m $\{\underline{s}_{ci}\}_{i=1}^m \sim sPPM$, but that's
not equivalent to $\rho_m \sim sPPM$, as for the ρ_m we
can transform to get the c_i .

$$\rho_m = \{s_{i1}, \dots, s_{im}\}$$

$\Leftrightarrow \begin{pmatrix} 1 \\ 2 \\ \vdots \\ n \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ \vdots \\ n \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ \vdots \\ n \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ \vdots \\ n \end{pmatrix}$
 $km = n$

car is worn $i:m$
| $c_i = \text{vehicle}(i: \text{wm } s_i \text{ for } j = e: \text{km})$
end

page 3
rec 22

On next, now we can understand better the model of an infinite most relevant prior, the π_{prior} .

Here the slow module the tree extension at count change

As the set $w=4 \rightarrow m$ the units, so, the model is recovered as $p_t = \{S_1^{(t)}, \dots, S_{k_t}^{(t)}\} \rightarrow \text{NICE}$ or each $S_t = (c_{1(t)}, \dots, c_{n(t)})$ as below.

So we want to study the IP of getting a certain combination of partitions. So, using a MC structure, we let our treelet be as never) structure, we let our treelet be as

$$\varphi(p_u, \dots, p_T) = P(P_u = p_u, \dots) = \varphi(p_1 | p_{u-1}) \cdot \dots \cdot \varphi(p_T | p_u) \cdot \varphi(p_u)$$

product rule

Then the we as you see the partition structure this

$$\Pr(p/M) = \varphi_{p_u|M}(p/M) = P(p_u = p/M) = \frac{M^u}{\prod_{i=1}^m (M+i-1)} \frac{\prod_{i=1}^u (1+i-1)!}{\text{tree constant}} \frac{\Pr(S_i)}{\Pr(S)}$$

told in the version 7

we call this law CRP(M).
But this law misses how we want to model the connection among c_{t-1} and c_t .
so we call it CRP(M)

So we introduce

$$\alpha_t = \{0 \mid \text{unit } i \text{ stays in the same cluster from } t-1 \text{ to } t\}$$

$\alpha_t \sim \text{Ber}(\alpha_t) / \begin{cases} \alpha_t = 0: p_t \perp c_{t-1} \\ \alpha_t = 1: p_t = c_{t-1} \end{cases}$

we get see as more $\alpha_t \sim \text{Ber}(\alpha_t)$

then as better uniform, creating we introduce
variables $\mathcal{I}_t = (\mathcal{I}_{1(t)}, \dots, \mathcal{I}_{m(t)})$ where $\mathcal{I}_{i(t)} \in \{0, 1, 2\}$, to
which values we set from the $\text{Ber}(\alpha_t)$ for the α_t
we introduced above.

We call this model tRPM(α, M), and later in what
we will end in the model, short st, we synthesize
at the p_t .

$$\begin{aligned} Y_{it} | \mu_t^*, \sigma_{it}^{2*}, c_t &\stackrel{\text{ind}}{\sim} N(\mu_{c_{it}t}^*, \sigma_{c_{it}t}^{2*}), i = 1, \dots, m \text{ and } t = 1, \dots, T, \\ (\mu_{jt}^*, \sigma_{jt}^*) | \theta_t, \tau_t^2 &\stackrel{\text{ind}}{\sim} N(\theta_t, \tau_t^2) \times \text{UN}(0, A_\sigma), j = 1, \dots, k_t, \\ (\theta_t, \tau_t) &\stackrel{\text{iid}}{\sim} N(\phi_0, \lambda^2) \times \text{UN}(0, A_\tau), t = 1, \dots, T, \\ (\phi_0, \lambda) &\sim N(m_0, s_0^2) \times \text{UN}(0, A_\lambda), \\ \checkmark \{c_t, \dots, c_T\} &\sim \text{tRPM}(\alpha, M), \text{ with } \alpha_t \stackrel{\text{iid}}{\sim} \text{Beta}(a_\alpha, b_\alpha), \end{aligned} \tag{5}$$

The effect of this \mathcal{I}_t on the law of p is just that it reduces the support, ie now

$$P(c_t = 1 | \mathcal{I}_t, p_{u-1}) = \{0 \mid \text{if } c_{t-1} \text{ and } c_t \text{ are not }\mathcal{I}_t\}$$

: now it gets more complex, less clear

Actually, the prior continues and becomes non parametric as it tells about certain laws from a "Dirichlet process (DP)", which means, as a "prior that outputs distributions".

So maybe its too complex.
So we can switch to our own models.

Or we), see there was a large separation mechanism
so we could give it to make set θ our own
new models.

$$\gamma(\omega) = \xi_{\omega} = \gamma_{\text{out}}^{\Phi}(\omega) + \xi_{\omega}^T \beta_{\text{out}} + \theta(\omega) + \xi_{\omega}^T \beta$$

} spatial effect
center + global specific

} covariate effect
cluster + global specific

page 3
sec 23

$$\xi_{\omega}(t) | \xi_{\omega}^{(t-4)}, \dots, \xi_t \stackrel{\perp}{\sim} N(\underbrace{\gamma_{\text{out}}^{\Phi}(t) + w_t \xi_{\omega}^{(t-4)}}_{\text{to introduce an AR(4) effect}}, \dots)$$

page 4
sec 40

also by drawing just \perp (we \perp)
we can use all \perp models as
mixture models

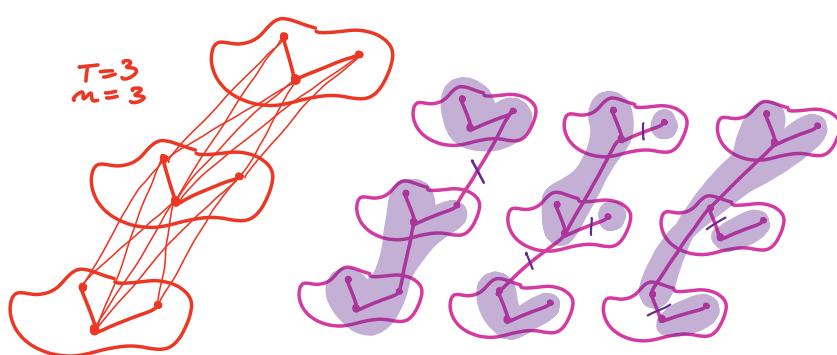
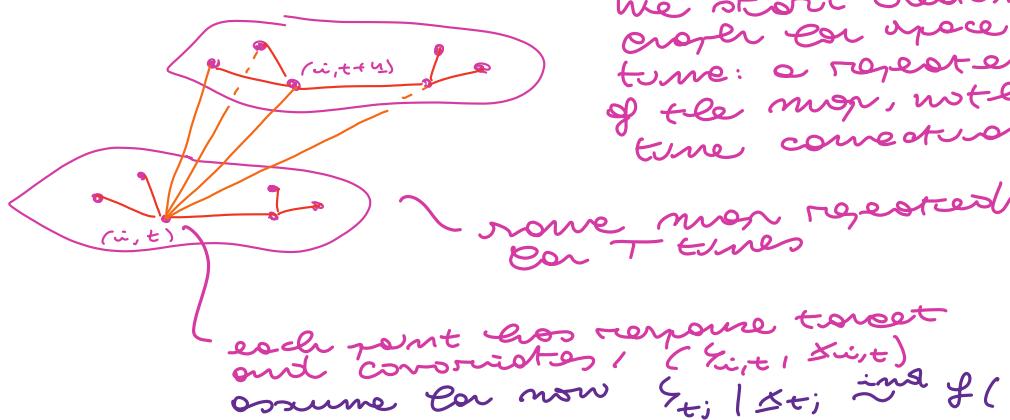
$$\left| \begin{array}{l} \xi_{\omega}(t) | w_t, f_{\omega} \sim \text{ind} \quad \text{and} \quad f_{\omega}(y) = \sum_{q=1}^{k(t)} w_{q,t} \cdot f_q(y) \\ \text{where } w_{q,t} = P(\omega_{q,t} = q | w) \\ \text{or} \\ \xi_{\omega}(t) | \xi_t, \dots \sim \text{ind} \dots \text{where } \xi_t \sim \dots \end{array} \right.$$

unstructured
TA γ

Model about the crash structure.

We start creating a
crash car space and
time: a repeated release
of the car, with also
time connections.

page 4



Where on page 4 we can use covariates in the rest of the car.
But for me I think it is not enough, as the car is one,
but covariates vary over time.

page 4