

**Politecnico di Milano**

School of Industrial and Information Engineering  
Master of Science in Mathematical Engineering

MASTER THESIS

The DRPM Strikes Back: ~~Improvements~~ *more flexibility for* on a  
Bayesian Spatio-Temporal Clustering Model

Advisor

**Prof. Alessandra Guglielmi**

Coadvisor

**Prof. Alessandro Carminati**

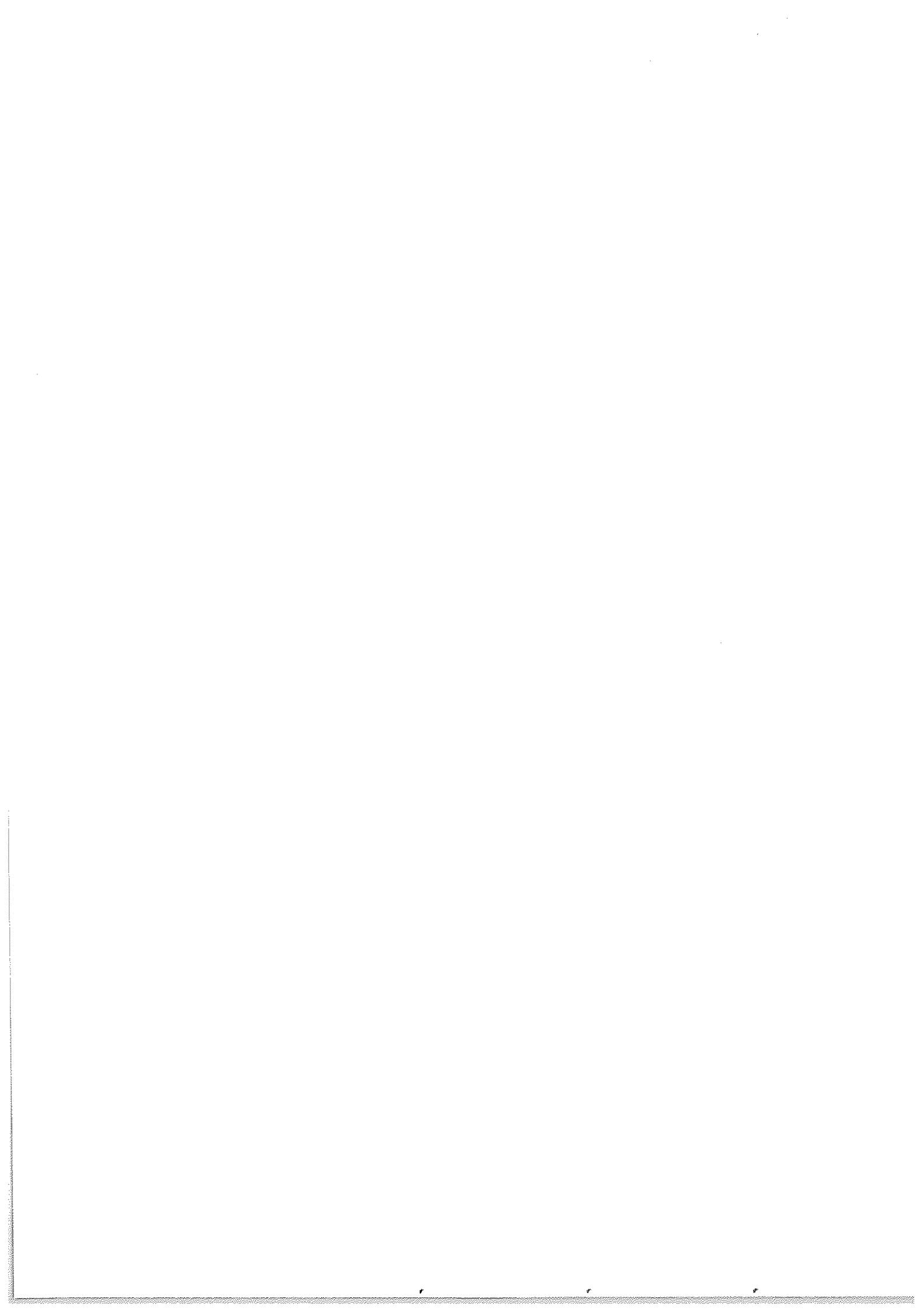
Candidate

**Federico Angelo Mor**

Matr. 221429



*to my cats Otto  
and La Micia*

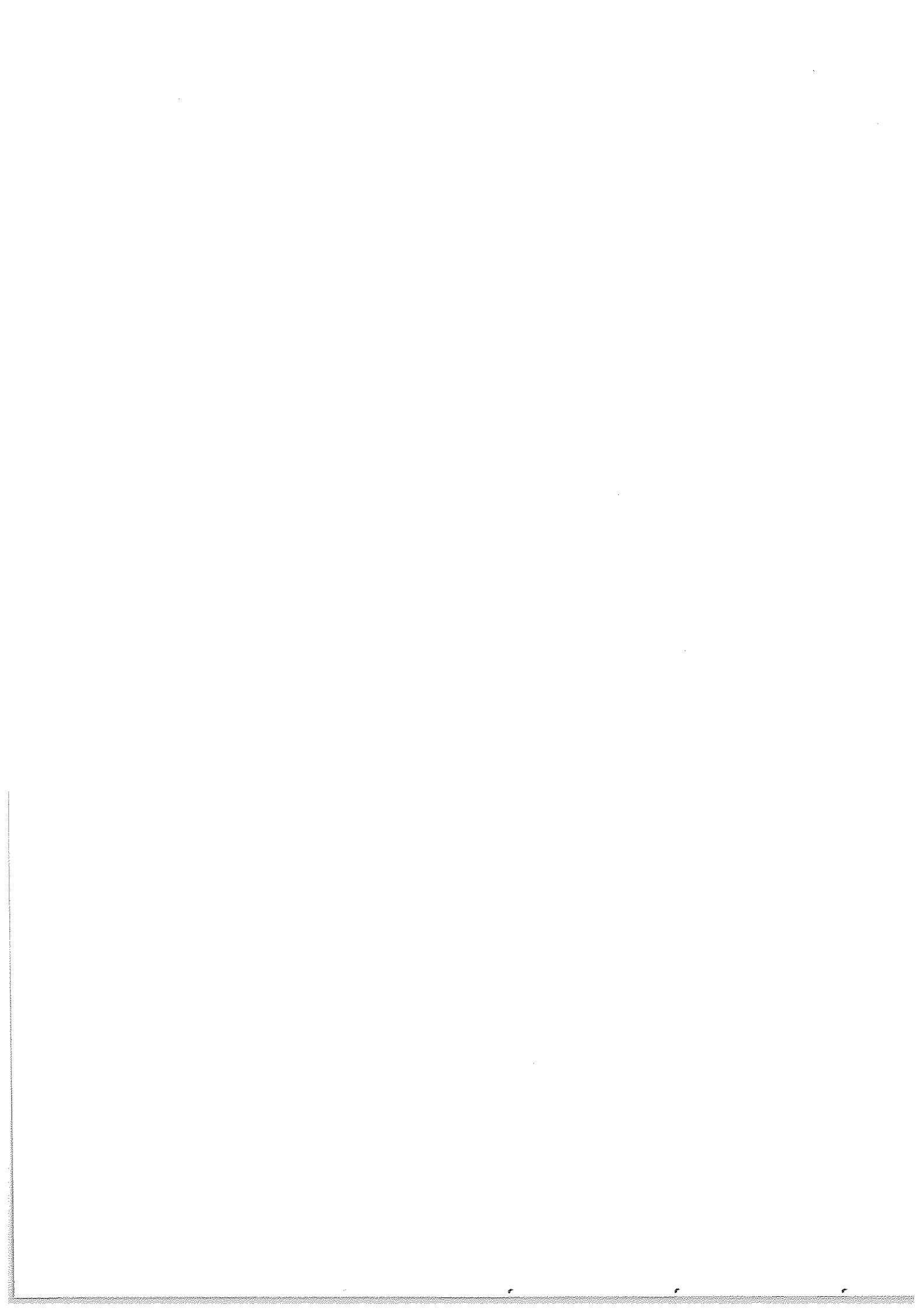


However, the current formulation of the model and the implementation of the associated MCMC lack

# Abstract

Clustering is a key technique for identifying patterns and structures in complex datasets, whose relevance is intensified in spatio-temporal contexts where observations are simultaneously influenced by multiple factors such as space, time, and covariates. To this end, the Dependent Random Partition Model (DRPM) is one of the most relevant Bayesian models due to its explicit consideration of temporal dependence in the partitions. However, the current implementation lacks of the inclusion of covariates, the handling of missing data, and the efficiency in execution times. Therefore, in this work we improve the original DRPM model by addressing those issues through updates on the model formulation and a brand new implementation in Julia. These advancements are then tested on synthetic and real-world datasets, including air quality data from the AgrImOnIA project in Lombardy, Italy.

KEYWORDS: Bayesian modelling, clustering, spatio-temporal data, computational statistics, MCMC, Julia



# Sommario

Vedi versione in INGLESE

Il clustering è una tecnica fondamentale per identificare strutture e pattern in dataset complessi, la cui importanza è intensificata nei contesti spazio-temporali in cui le osservazioni sono influenzate simultaneamente da molteplici fattori come spazio, tempo e covariate. In tal senso, il modello DRPM (Dependent Random Partition Model, modello per partizioni aleatorie dipendenti) è uno dei modelli bayesiani più rilevanti in quanto tiene conto in modo esplicito della dipendenza temporale delle partizioni. Tuttavia, l'attuale implementazione manca dell'inclusione di covariate, della gestione dei dati mancanti, e di efficienza nei tempi di esecuzione. In questo lavoro abbiamo quindi migliorato l'originale modello DRPM affrontando tali problemi tramite aggiornamenti sulla formulazione del modello e una fiammante implementazione in Julia. Questi sviluppi sono stati poi testati su dataset sintetici e reali, compresi i dati sulla qualità dell'aria in Lombardia del progetto AgrImOnIA.

PAROLE CHIAVE: modello bayesiano, clustering, dati spazio-temporali, statistica computazionale



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Description of the model</b>	<b>3</b>
2.1	Update rules derivation . . . . .	7
2.2	Spatial cohesions analysis . . . . .	10
2.3	Covariates similarities analysis . . . . .	13
<b>3</b>	<b>Implementation and optimizations</b>	<b>21</b>
3.1	Optimizations . . . . .	22
3.1.1	Optimizing spatial cohesions . . . . .	23
3.1.2	Optimizing covariates similarities . . . . .	26
<b>4</b>	<b>Testing</b>	<b>29</b>
4.1	Assessing the equivalence of the models . . . . .	29
4.1.1	Target variable only . . . . .	30
4.1.2	Target variable plus space . . . . .	31
4.2	Performance with missing values . . . . .	38
4.2.1	Target variable only (NA case) . . . . .	38
4.2.2	Target variable plus space (NA case) . . . . .	40
4.3	Effects of the covariates . . . . .	41
4.3.1	Covariates in the likelihood . . . . .	43
4.3.2	Covariates in the clustering . . . . .	48
4.4	Scaling performances . . . . .	52
<b>5</b>	<b>Conclusion</b>	<b>59</b>
<b>A</b>	<b>Theoretical details</b>	<b>61</b>
A.1	Extended computations of the full conditionals . . . . .	61

<b>B Computational details</b>	<b>71</b>
B.1 Fitting algorithm code . . . . .	71
B.2 Interface . . . . .	86
<b>C Fits interpretation</b>	<b>91</b>
<b>Bibliography</b>	<b>93</b>

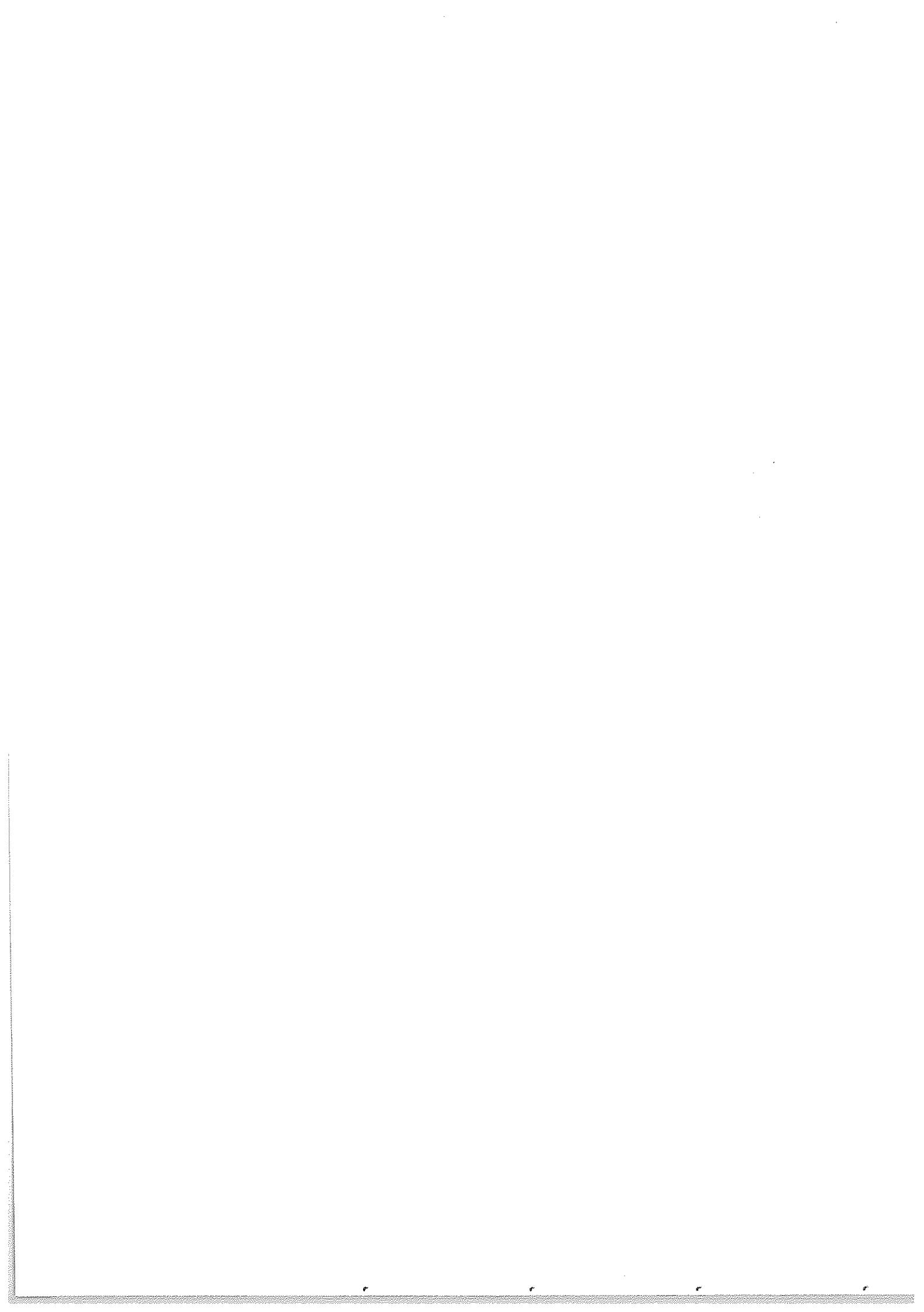
# List of Figures

2.1	Updated DRPM model graph . . . . .	6
2.2	Partition considered for the cohesion analysis . . . . .	12
2.3	Cohesions 1, 2, and 3 illustration . . . . .	14
2.4	Cohesions 4, 5, and 6 illustration . . . . .	15
2.5	Partition considered for the similarity analysis . . . . .	16
2.6	Similarities 1, 2 and 3 illustration . . . . .	18
2.7	Similarity 4 illustration . . . . .	19
3.1	Flame graph of a test fit . . . . .	22
3.2	Cohesions 3 and 4 implementation comparison . . . . .	25
3.3	Similarity 4 annotations comparison . . . . .	27
4.1	Lagged ARI values of CDRPM and JDRPM fits, target values only	30
4.2	Clusters produced by JDRPM and CDRPM fits, target values only	31
4.3	Visual representation of the clusters of JDRPM and CDRPM fits, target values only . . . . .	32
4.4	Generated and fitted values of JDRPM and CDRPM fits, target values only . . . . .	33
4.5	Comparison of the two possible mean centering methods on the target variable . . . . .	34
4.6	Lagged ARI values of CDRPM and JDRPM fits, target plus space values . . . . .	35
4.7	Target and fitted values of JDRPM and CDRPM fits, target plus space values . . . . .	36
4.8	Clusters generated by CDRPM and JDRPM fits, target plus space values . . . . .	37
4.9	Fitted values of JDRPM fit, target values only, NA dataset . . . . .	38
4.10	Clusters produced by JDRPM fits, target values only, full vs NA dataset . . . . .	39

4.11 Lagged ARI values of JDRPM fits, target values only, full vs NA dataset . . . . .	39
4.12 Visual representation of the clusters of JDRPM fit, target values only, NA dataset . . . . .	40
4.13 Lagged ARI values of JDRPM fits, target plus space values, full vs NA dataset . . . . .	41
4.14 Fitted values of JDRPM fit, target plus space values, NA dataset . . . . .	41
4.15 Comparison of the two possible mean centering methods on a covariate . . . . .	42
4.18 Lagged ARI values of JDRPM fit, target plus space values, full vs NA dataset, with covariates in the likelihood . . . . .	43
4.16 Regression vector of the fit with multiple covariates in the likelihood, full dataset . . . . .	44
4.17 Regression vector of the fit with multiple covariates in the likelihood, NA dataset . . . . .	45
4.19 Clusters generated by JDRPM fits, target plus space values, full vs NA dataset, with covariates in the likelihood . . . . .	46
4.20 Target and fitted values of JDRPM fits, target plus space values, NA dataset, with covariates in the likelihood . . . . .	47
4.21 Trace plot of the fitted values for a fit with covariates in the likelihood . . . . .	47
4.22 Variables used for the fit tests with covariates in the clustering . . . . .	49
4.23 Comparison of the clusterings provided by the base fits plus JDRPM with covariates in the clustering . . . . .	51
4.24 Clusters generated by JDRPM fit, target plus space values, with covariates in the clustering process . . . . .	52
4.25 CDRPM standard fit, clusters distribution with respect to the wind speed covariate . . . . .	53
4.26 JDRPM fit with covariates, clusters distribution with respect to the wind speed covariate . . . . .	53
4.27 JDRPM standard fit, clusters distribution with respect to the wind speed covariate . . . . .	54
4.28 Execution times of JDRPM and CDRPM fits, target values only . . . . .	55
4.29 Execution times of JDRPM and CDRPM fits, target plus space values . . . . .	55
4.30 Execution times of JDRPM fits, target plus space plus covariates . . . . .	56
4.31 Visual representation of all fitting performances . . . . .	57

# List of Tables

4.1	Accuracy metrics of CDRPM and JDRPM fits, target values only . . . . .	30
4.2	Accuracy metrics of CDRPM and JDRPM fits, target plus space values . . . . .	35
4.3	Accuracy metrics of JDRPM fits, target values only, full vs NA dataset . . . . .	39
4.4	Accuracy metrics of JDRPM fits, target plus space values, full vs NA dataset . . . . .	40
4.5	Accuracy metrics of JDRPM fits, target plus space values, full vs NA dataset, with vs without covariates in the likelihood . . . . .	43
4.6	Accuracy metrics of CDRPM and JDRPM fits, target plus space values, with vs without covariates in the clustering process . . . . .	50



Per prima cosa vuol dire seguire un approccio model-based al clustering (per es. significa usare le misure per come le unità dei dati, oppure n'èice dare le loro migliore caratteristiche al parametru elettrico),

poi cose vuol dire usare un modello bayesiano per fare clustering (ad esempio i parametri elettrici (incluso le pertinenze degli individui)). Dire che è necessario

## Chapter 1

disegnare metodi MCMC per approssimare la posterior (su cui lavoreremo tutta l'inferenza); tali metodi sono molti corssi del

### Introduction

problema numerico in generale. Per cominciare a dire del Bayesian model per il clustering di dati spazio-tempo

bayesian = massiccio in ITALIANO; Bayesian = massiccio in INGLESE

Clustering has always been a powerful tool to identify structures and patterns in data especially in contexts where relationships between the observations are complex, e.g. when the target variable is affected by many factors simultaneously. For this reason, clustering techniques saw a continuous increase in popularity in a variety of scientific fields, including social sciences, climate and environmental analysis, economics, and healthcare. The importance of clustering becomes even more noticeable when working on spatio-temporal datasets, in which observations are collected over time and across different spatial locations, possibly concealing trends behind both information levels. This type of data, in fact, is inherently complex due to this dependence and interaction between spatial and temporal dimensions; a complexity that is further increased if covariates are also available. For that reason, an effective analysis of such data demands models that can account for this dependence while also providing efficient implementations to be possibly applied on large scale datasets, which are commonly accessible in this context.

Recently, the use of Bayesian ~~models~~ to perform clustering has gained some attention, particularly in this field of spatio-temporal datasets. Bayesian clustering, in fact, allows to incorporate prior information into the model enhancing the flexibility and interpretability of the results with respect, for example, to more traditional frequentist approaches. Throughout the years, several models have been developed, but one of the most relevant to this end is the Dependent Random Partition Model (DRPM), which stands out for being able to handle explicitly the temporal dependence of partitions into the model formulation, while also possibly accounting for the spatial information. However, the current DRPM implementation, written in C and available through an R interface, lacks some relevant utilities such as the inclusion of covariates, which could further improve the generation and informativity of the clusters, the handling of missing data, and an efficient implementation, which would speed up the model fitting to e.g. run multiple chains in parallel or be more easily applied on large scale datasets.

In this work, we aim to address these three issues by enlarging the original model, that is, preserving the primary idea of the formulation but making it richer through the insertion of new components. We will show how this updated model can perform better than the original one, under the same testing conditions, and can also

*che o più ??*

improve the clustering accuracy and interpretation through the inclusion of covariates. All this while also providing faster execution times. In fact, implementing the model using the Julia language, rather than C, we took advantage of its high-performance capabilities and well-equipped statistical ecosystem. Our comparison will focus on both synthetic datasets and real-world applications, with the latter involving air quality measurements from the AgriMoIA dataset, a comprehensive record of air pollutant levels and other environmental variables measured across the Lombardy region of Italy.

Chapter 1 will briefly review the literature about bayesian clustering models and then dive deeply into the analysis and description of DRPM, and of our updated version, about their core aspects of sampling algorithm, spatial cohesions, and covariates similarities. *non a se encare cose nuove!*

Chapter 2 will provide some insights about the computational aspects of the model implementation, motivating the choice of the Julia language and reporting some optimization possibilities emerged when developing the algorithm.

Chapter 3 will be devoted to test and compare the original DRPM formulation and implementation to our updated version. We will check if they perform similarly, at a common testing level, and assess the performances of our model when we employ the new updates, i.e. the handling of missing data and the insertion of covariates at clustering and likelihood levels. An analysis about expected execution times with respect to the size of the dataset will also be provided.

Finally, in Chapter 4, we will briefly review the benefits and drawbacks that this work revealed and suggest possible further improvements or development paths.

??

? are test? *il nostro algoritmo*  
 ? cose confrontiamo: *il nostro modello algoritmo*  
     e *il DRPM modello e algoritmo originale*

*inglese!*

*generalization*

*categoria*

*for spatio-temporal data*

*MCMC*

*inglese!*

*? qualche EXPECTED*

Bayesian models for clustering are grouped in two classes of models.

The first one is based on extremes data [or random effects] and

points

Def di  
clustering  
su WIKIPEDIA

is distributed from a mixture density. The clustering labels, i.e. labels which identify which component of the mixture each point is associated to, or identify the clustering of the individuals units.

Chapter 2 The second class, instead, assumes specifies the conditional distribution of the data points given a realization of the partition of all the units, and a prior is assigned to this partition.

Mettere in RIF BIBLIOGRAFICI

## Description of the model

Ambrosio-Villejos et Walker (2015)

"Come on, gentlemen, why shouldn't we get rid of all this calm reasonableness with one good kick, just so as to send all these logarithms to the devil and be able to live our own lives at our own sweet will?"

— Fëodor Dostoevskij, Notes from the Underground

In the Bayesian framework, clustering is possible by employing a random probability measure of discrete type that induces a distribution over random partitions. This discreteness is obtained with the Dirichlet Process (DP), which several clustering models implement either through the stick-breaking representation [Bar+12] [AW15] [GMR16] [Jo+16] [KG18] [DK18] [De +19] or through the Pólya urn scheme [Car+17]. However, these classical bayesian methods rely on modelling the dependence in the random partitions by modelling the dependence inside the random probability measures, i.e. on the parameters which underlie those DP representations rather than to the clusters themselves. This approach is therefore kind of a "step back" from the main object of interest, the clusters, which are then only induced by the random partition model. As a consequence, there is no guarantee that the correlation that appears in the parameters would subsequently reflect into correlation among the partitions, often producing counterintuitive behaviours in the results. The Dependent Random Partition Model (DRPM) [PQD22], on the other hand, models directly the sequence of partitions, thus providing a more reasonable, accurate, and interpretable temporal evolution of the clusters.

quelle correlation? Quelle temporal? Me encore non he un solo model

Before diving into the model description, we define some notation conventions. We setup in a spatio-temporal context with  $i = 1, \dots, n$  and  $t = 1, \dots, T$  being the indexes for units and time instants. We will denote with  $\rho_t = \{S_{1t}, \dots, S_{kt}\}$  the partition at time  $t$ , of the  $n$  experimental units, composed by  $k_t$  cluster. Another possible representation of the partition is through cluster membership labels  $c_t = \{c_{1t}, \dots, c_{nt}\}$ , where  $c_{it} = j$  if unit  $i$  belongs to cluster  $S_{jt}$ . Finally, we will denote with a  $*$  superscript all the variables or quantities which are cluster-specific.

To implement dependence in the partitions one could simply propose a joint probability model for  $(\rho_1, \dots, \rho_T)$ , denoted as  $P(\rho_1, \dots, \rho_T)$ , where each  $\rho_t$  is set to be possibly affected by all the other partitions. This principle, however,

needs specify NO former context  
clustered at all time points  
 $t = 1, 2, \dots, T$ . The units are represented as  $i = 1, 2, \dots, n$ .

??

To focus to have a require: è meglio spiegare le  
prior (dove le sue def complete) e poi eventualmente  
spiegare come significa avere delle hewe

Page et al (2022)

Chapter 2. Description of the model

too complex and general to be modelled efficiently; therefore the DRPM authors limited this temporal connection to a first-order Markov-chain structure, where the conditional distribution of  $\rho_t$  given all the predecessors  $\rho_{t-1}, \rho_{t-2}, \dots, \rho_1$  actually depends only on  $\rho_{t-1}$ . This brings the random partition model to the form

$$P(\rho_1, \dots, \rho_T) = P(\rho_T | \rho_{T-1}) \cdots P(\rho_2 | \rho_1) P(\rho_1) \quad (2.1)$$

To explicitly manage the relation between  $\rho_t$  and  $\rho_{t-1}$  some auxiliary variables are introduced. The idea is that if two partitions are highly time-dependent, few changes will occur between them. In turn, partitions which are quite independent will possibly exhibit very different configurations. To express this fixity or flexibility concept, for each unit  $i = 1, \dots, n$  the following variable is introduced<sup>1</sup> *no foot note*

$$\gamma_{it} = \begin{cases} 1 & \text{if unit } i \text{ is not reallocated when moving from time } t-1 \text{ to } t \\ 0 & \text{otherwise (i.e. unit } i \text{ is reallocated)} \end{cases} \quad (2.2)$$

By construction, we set  $\gamma_{i1} = 0$  for all  $i$ , meaning that at the first time instant all units will get reallocated since they have no partition to which they could be possibly fixed at. Regarding their modelling, the authors proposed  $\gamma_{it} \stackrel{\text{ind}}{\sim} \text{Ber}(\alpha_t)$  with  $\alpha_t \in [0, 1]$  behaving as a temporal dependence parameter. At the two extremes,  $\alpha_t = 1$  will denote perfect temporal dependence, with  $\rho_t = \rho_{t-1}$ , while  $\alpha_t = 0$  will imply full independence of  $\rho_t$  from  $\rho_{t-1}$ . For the sake of clarity, the vector  $\gamma_t = (\gamma_{1t}, \dots, \gamma_{nt})$  is created, and the augmented joint model becomes in the form

$$P(\gamma_1, \rho_1, \dots, \gamma_T, \rho_T) = P(\rho_T | \gamma_T, \rho_{T-1}) P(\gamma_T) \cdots P(\rho_2 | \gamma_2, \rho_1) P(\gamma_2) P(\rho_1) \quad (2.3)$$

inglese!

The insertion of these additional variables makes the model very powerful in describing the temporal dependence of the partitions but slightly hinders the design of the sampling algorithm. To outline it, we firstly need a ??

**Definition 2.1** (compatibility). Two partitions  $\rho_t$  and  $\rho_{t-1}$  are *compatible* with respect to  $\gamma_t$  if  $\rho_t$  can be obtained from  $\rho_{t-1}$  by reallocating items as indicated by  $\gamma_t$ ; i.e. only moving the units  $i$  with  $\gamma_{it} = 0$ .

To perform this compatibility check, it is enough to ensure that the reduced partitions from  $\rho_t$  and  $\rho_{t-1}$  are the same, with reduced meaning their restriction to the units which cannot move. Indeed, if those fixed units are clustered in the same ways, then surely the free-movers from  $\rho_t$  can be set to match the labels assigned by partition  $\rho_{t-1}$ . Denoting as  $\mathfrak{R}_t = \{i : \gamma_{it} = 1\}$  the set of fixed units at time  $t$ , this check translates into asking that  $\rho_t^{\mathfrak{R}_t} = \rho_{t-1}^{\mathfrak{R}_t}$ . *che buono!!*

The sampling algorithm requires that when we are drawing the new samples for the  $\gamma_{ita}$ , or also for the cluster labels  $c_{it}$ , we firstly need to check if those draws can actually be valid, i.e. if they would keep compatible and coherent all the partitions and parameters involved. For example, when updating  $\gamma_{it}$  during each iteration  $d$  of the algorithm, the only case which can raise problems is when we pass from

<sup>1</sup>a quick way to remind this convention is thinking that  $\gamma_{it}$  answers to the question "do I stay fixed?" asked from unit's  $i$  perspective.

quale?? Dimostro, meglio che lei enzigni prima tutto il modello,  
comprese le prior per  $(\gamma_1, \gamma_2, \dots, \gamma_T, \rho_T)$  e poi spieghi queste cose!

??

inglese  
inglese  
non si segue!

$\gamma_{it}^{(d-1)} = 0$  to  $\gamma_{it}^{(d)} = 1$ . This step corresponds to the case in which a unit  $i$  that was initially (i.e. according to the previous iteration parameters) free to be reassigned is now instead deemed to stay fixed in her cluster. However, this change may not align to the current sampled values of the partitions  $\rho_{t-1}^{(d)}$  and  $\rho_t^{(d-1)}$ . Therefore, compatibility between their reductions to the units in the set  $\mathfrak{R}_t \cup \{i\}$  needs to be checked and, if this check fails, the tentative update  $\gamma_{it}^{(d-1)} = 0 \rightarrow \gamma_{it}^{(d)} = 1$  is not allowed to be performed and we force  $\gamma_{it}^{(d)} = 0$  in the sampling algorithm. Similar checks are conducted when  $\rho_t$  is updated. In this step, only the units that can actually move, i.e. that have  $\gamma_{it} = 0$ , are updated, and therefore there are no compatibility problems between  $\rho_{t-1}$  and  $\rho_t$ . However, since the update of  $\gamma_{it}$  occurs before the one of the partition, compatibility needs to be checked between  $\rho_t$  and  $\rho_{t+1}$ .

dopo  
la decisione  
dell'algoritmo

In any case, once the partition model is specified, there is great flexibility in how to setup the rest of the hierarchical model. To allow temporal dependence to propagate through the model, an autoregressive AR(1) component is also added to the formulation of the model (but only optionally in the implementation), both at the data and cluster-specific parameters level. All this led the authors to the following complete model

$$\begin{aligned}
 Y_{it}|Y_{it-1}, \mu_t^*, \sigma_t^{2*}, \eta, c_t &\stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_{c_{it}}^* + \eta_{1i} Y_{it-1}, \sigma_{c_{it}}^{2*}(1 - \eta_{1i}^2)) & i = ? \quad t = \dots? \\
 Y_{i1} &\stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_{c_{i1}}^*, \sigma_{c_{i1}}^{2*}) \\
 \xi_i = \text{Logit}(\frac{1}{2}(\eta_{1i} + 1)) &\stackrel{\text{ind}}{\sim} \text{Laplace}(a, b) \\
 (\mu_{jt}^*, \sigma_{jt}^*) &\stackrel{\text{ind}}{\sim} \mathcal{N}(\vartheta_t, \tau_t^2) \times \mathcal{U}(0, A_\sigma) \\
 \vartheta_t | \vartheta_{t-1} &\stackrel{\text{ind}}{\sim} \mathcal{N}((1 - \varphi_1)\varphi_0 + \varphi_1 \vartheta_{t-1}, \lambda^2(1 - \varphi_1^2)) \\
 (\varphi_1, \tau_t) &\stackrel{\text{iid}}{\sim} \mathcal{N}(\varphi_0, \lambda^2) \times \mathcal{U}(0, A_\tau) \\
 (\varphi_0, \varphi_1, \lambda) &\sim \mathcal{N}(m_0, s_0^2) \times \mathcal{U}(-1, 1) \times \mathcal{U}(0, A_\lambda) \\
 \{c_t, \dots, c_T\} &\sim \text{tRPM}(\alpha, M) \text{ with } \alpha_t \stackrel{\text{iid}}{\sim} \text{Beta}(a_\alpha, b_\alpha) \tag{2.4}
 \end{aligned}$$

where tRPM represents the temporal random partition model (2.3).

Moving towards our update, we decided to refine some parts of that formulation. Regarding the variances  $\sigma_{jt}^{2*}$ ,  $\tau_t^2$ , and  $\lambda^2$ , we chose to model them through an inverse gamma distribution rather than the uniform employed originally. This is indeed a more sophisticated choice, since the tuning of the parameters of an  $\text{invGamma}(a, b)$  is a bit more difficult than simply setting the bounds of a  $\mathcal{U}(l, u)$ , but should guarantee a better mixing in the chain. In fact, the invGamma distributions recovers conjugacy in the model, thanks to the normal law assigned to the other parameters, allowing the update step of the variances to be performed through the analytically exact Gibbs sampler rather than the acceptance-rejection method of Metropolis algorithm. Finally, to improve the accuracy in fitting the target values, we added a regression parameter  $\beta_t$  in the likelihood. We decided to make it only time-dependent, and not also unit-dependent, to lighten the already quite-heavy formulation.

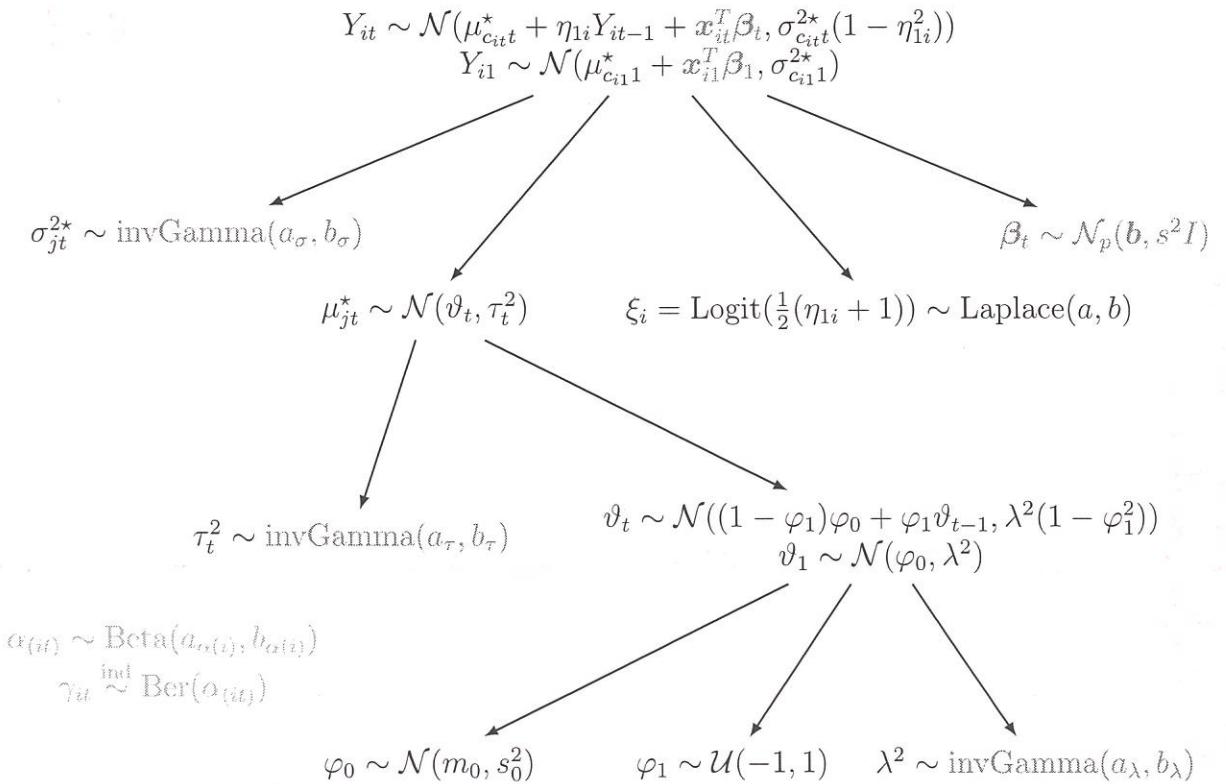
The final updated model is now proposed, with highlighted in dark red the

Facciamo le denunce anche [con formula]  
del nostro modello, e poi fece i commenti! Sarà  
che non ci capisce niente!

changes and insertions that we made.

$$\begin{aligned}
 Y_{it} | Y_{it-1}, \mu_t^*, \sigma_{jt}^{2*}, \eta, c_t &\stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_{c_{it}}^* + \eta_{1i} Y_{it-1} + x_{it}^T \beta_t, \sigma_{c_{it}}^{2*}(1 - \eta_{1i}^2)) \\
 Y_{i1} &\stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_{c_{i1}}^* + x_{i1}^T \beta_1, \sigma_{c_{i1}}^{2*}) \\
 \beta_t &\stackrel{\text{ind}}{\sim} \mathcal{N}_p(\mathbf{b}, s^2 I) \\
 \xi_i = \text{Logit}(\frac{1}{2}(\eta_{1i} + 1)) &\stackrel{\text{ind}}{\sim} \text{Laplace}(a, b) \\
 (\mu_{jt}^*, \sigma_{jt}^{2*}) &\stackrel{\text{ind}}{\sim} \mathcal{N}(\vartheta_t, \tau_t^2) \times \text{invGamma}(a_\sigma, b_\sigma) \\
 \vartheta_t | \vartheta_{t-1} &\stackrel{\text{ind}}{\sim} \mathcal{N}((1 - \varphi_1)\varphi_0 + \varphi_1 \vartheta_{t-1}, \lambda^2(1 - \varphi_1^2)) \\
 (\vartheta_1, \tau_t^2) &\stackrel{\text{iid}}{\sim} \mathcal{N}(\varphi_0, \lambda^2) \times \text{invGamma}(a_\tau, b_\tau) \\
 (\varphi_0, \varphi_1, \lambda^2) &\sim \mathcal{N}(m_0, s_0^2) \times \mathcal{U}(-1, 1) \times \text{invGamma}(a_\lambda, b_\lambda) \\
 \{c_t, \dots, c_T\} &\sim \text{tRPM}(\alpha, M) \text{ with } \alpha_t \stackrel{\text{iid}}{\sim} \text{Beta}(a_\alpha, b_\alpha)
 \end{aligned} \tag{2.5}$$

A visual representation of this new version of the DRPM model is also present in Figure 2.1, to more clearly appreciate the hierarchical structure and the relations among the parameters.



**Figure 2.1:** Graph visualization of the DRPM model, with highlighted in dark red the changes that we made to the original formulation and in gray the internal variables of the model.

In the course of this work, for the sake of clarity, we will refer to CDRPM for the original model formulation of [PQD22] and to JDRPM for our updated version.

We will now dive more deeply into the characteristics of the models by deriving the update rules for the parameters which will be used to implement the MCMC fit-

we will use MCMC.

ting algorithm and, subsequently, inspecting the behaviours of the spatial cohesions and covariates similarities.

## 2.1 Update rules derivation

?? The DRPM Gibbs sampler  
o è quello di PV no?!

Io scrivo  
tutte le cose

o almeno  
descrivrei  
per me un modello

di DRPM, poi  
l'avevate già fatto,  
poi il resto  
modello, poi  
algoritmo mat  
etc

We briefly report the full conditionals derivation for the parameters which had a conjugacy in the model (for the full computations see Appendix A). The other variables not included here, namely  $\eta_{1i}$  and  $\varphi_1$ , involved instead the classical Metropolis update.

Indicare gli step nuovi o DIVERSI  
rispetto a Pege et al.

- update  $\sigma_{jt}^{2*}$

for  $t = 1$ :  $f(\sigma_{jt}^{2*} | -) \propto$  kernel of a  $\text{invGamma}(a_{\sigma(\text{post})}, b_{\sigma(\text{post})})$  with

$$a_{\tau(\text{post})} = a_\sigma + \frac{|S_{jt}|}{2} \quad b_{\tau(\text{post})} = b_\sigma + \frac{1}{2} \sum_{i \in S_{jt}} (Y_{it} - \mu_{jt}^* - \mathbf{x}_{it}^T \boldsymbol{\beta}_t)^2$$

for  $t > 1$ :  $f(\sigma_{jt}^{2*} | -) \propto$  kernel of a  $\text{invGamma}(a_{\sigma(\text{post})}, b_{\sigma(\text{post})})$  with

$$a_{\tau(\text{post})} = a_\sigma + \frac{|S_{jt}|}{2} \quad b_{\tau(\text{post})} = b_\sigma + \frac{1}{2} \sum_{i \in S_{jt}} (Y_{it} - \mu_{jt}^* - \eta_{1i} Y_{it-1} - \mathbf{x}_{it}^T \boldsymbol{\beta}_t)^2$$

(2.6)

- update  $\mu_{jt}^*$

for  $t = 1$ :  $f(\mu_{jt}^* | -) \propto$  kernel of a  $\mathcal{N}(\mu_{\mu_{jt}^*(\text{post})}, \sigma_{\mu_{jt}^*(\text{post})}^2)$  with

$$\sigma_{\mu_{jt}^*(\text{post})}^2 = \frac{1}{\frac{1}{\tau_t^2} + \frac{|S_{jt}|}{\sigma_{jt}^{2*}}} \quad \mu_{\mu_{jt}^*(\text{post})} = \sigma_{\mu_{jt}^*(\text{post})}^2 \left( \frac{\vartheta_t}{\tau_t^2} + \frac{\sum_{i \in S_{jt}} (Y_{i1} - \mathbf{x}_{it}^T \boldsymbol{\beta}_t)}{\sigma_{jt}^{2*}} \right)$$

for  $t > 1$ :  $f(\mu_{jt}^* | -) \propto$  kernel of a  $\mathcal{N}(\mu_{\mu_{jt}^*(\text{post})}, \sigma_{\mu_{jt}^*(\text{post})}^2)$  with

$$\sigma_{\mu_{jt}^*(\text{post})}^2 = \frac{1}{\frac{1}{\tau_t^2} + \frac{\sum_{i \in S_{jt}} \frac{1}{1-\eta_{1i}^2}}{\sigma_{jt}^{2*}}} \quad \mu_{\mu_{jt}^*(\text{post})} = \sigma_{\mu_{jt}^*(\text{post})}^2 \left( \frac{\vartheta_t}{\tau_t^2} + \frac{\sum_{i \in S_{jt}} \frac{Y_{it}-\eta_{1i}Y_{i,t-1}-\mathbf{x}_{it}^T \boldsymbol{\beta}_t}{1-\eta_{1i}^2}}{\sigma_{jt}^{2*}} \right)$$

(2.7)

- update  $\boldsymbol{\beta}_t$

for  $t = 1$ :  $f(\boldsymbol{\beta}_t | -) \propto$  kernel of a  $\mathcal{N}(\mathbf{b}_{(\text{post})}, A_{(\text{post})})$  with

$$A_{(\text{post})} = \left( \frac{1}{s^2} I + \sum_{i=1}^n \frac{\mathbf{x}_{it} \mathbf{x}_{it}^T}{\sigma_{c_{it}^*}^2} \right)^{-1} \quad \mathbf{b}_{(\text{post})} = A_{(\text{post})} \left( \frac{\mathbf{b}}{s^2} + \sum_{i=1}^n \frac{(Y_{it} - \mu_{c_{it}^*}^*) \mathbf{x}_{it}}{\sigma_{c_{it}^*}^2} \right)$$

i.e.  $f(\boldsymbol{\beta}_t | -) \propto$  kernel of a  $\mathcal{N}\text{Canon}(\mathbf{h}_{(\text{post})}, J_{(\text{post})})$  with

$$\mathbf{h}_{(\text{post})} = \left( \frac{\mathbf{b}}{s^2} + \sum_{i=1}^n \frac{(Y_{it} - \mu_{c_{it}^*}^*) \mathbf{x}_{it}}{\sigma_{c_{it}^*}^2} \right) \quad J_{(\text{post})} = \left( \frac{1}{s^2} I + \sum_{i=1}^n \frac{\mathbf{x}_{it} \mathbf{x}_{it}^T}{\sigma_{c_{it}^*}^2} \right)$$

for  $t > 1$ :  $f(\boldsymbol{\beta}_t | -) \propto$  kernel of a  $\mathcal{N}(\mathbf{b}_{(\text{post})}, A_{(\text{post})})$  with

$$A_{(\text{post})} = \left( \frac{1}{s^2} I + \sum_{i=1}^n \frac{\mathbf{x}_{it} \mathbf{x}_{it}^T}{\sigma_{c_{it}t}^{2*}} \right)^{-1} \quad \mathbf{b}_{(\text{post})} = A_{(\text{post})} \left( \frac{\mathbf{b}}{s^2} + \sum_{i=1}^n \frac{(Y_{it} - \mu_{c_{it}t}^* - \eta_{1i} Y_{it-1}) \mathbf{x}_{it}}{\sigma_{c_{it}t}^{2*}} \right)$$

i.e.  $f(\beta_t | -) \propto \text{kernel of a } \mathcal{N}\text{Canon}(\mathbf{h}_{(\text{post})}, J_{(\text{post})})$  with

$$\mathbf{h}_{(\text{post})} = \left( \frac{\mathbf{b}}{s^2} + \sum_{i=1}^n \frac{(Y_{it} - \mu_{c_{it}t}^* - \eta_{1i} Y_{it-1}) \mathbf{x}_{it}}{\sigma_{c_{it}t}^{2*}} \right) \quad J_{(\text{post})} = \left( \frac{1}{s^2} I + \sum_{i=1}^n \frac{\mathbf{x}_{it} \mathbf{x}_{it}^T}{\sigma_{c_{it}t}^{2*}} \right) \quad (2.8)$$

Here  $\mathcal{N}\text{Canon}(\mathbf{h}, J)$  is the canonical formulation of the  $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$ , with  $\mathbf{h} = \Sigma^{-1} \boldsymbol{\mu}$  and  $J = \Sigma^{-1}$ . This other distribution facilitates the sampling, since these full conditional computations allow to derive directly the parameters of the canonical one, e.g. the inverse of the variance matrix rather than the variance matrix itself; and therefore sampling through it does not require any inversion of matrices which would produce more computational load, numerical instabilities, and loss of accuracy. As a consequence, in Julia we can write `rand(MvNormalCanon(h_star, J_star))` rather than the riskier one `rand(MvNormal(inv(J_star)*h_star, inv(J_star)))`; which apart from the previously mentioned disadvantages would be a statistically equivalent form.

- update  $\tau_t^2$

$f(\tau_t^2 | -) \propto \text{kernel of a } \text{invGamma}(a_{\tau(\text{post})}, b_{\tau(\text{post})})$  with

$$a_{\tau(\text{post})} = \frac{k_t}{2} + a_\tau \quad b_{\tau(\text{post})} = \frac{\sum_{j=1}^{k_t} (\mu_{jt}^* - \vartheta_t)^2}{2} + b_\tau \quad (2.9)$$

- update  $\vartheta_t$

for  $t = T$ :  $f(\vartheta_t | -) \propto \text{kernel of a } \mathcal{N}(\mu_{\vartheta_t(\text{post})}, \sigma_{\vartheta_t(\text{post})}^2)$  with

$$\sigma_{\vartheta_t(\text{post})}^2 = \frac{1}{\frac{1}{\lambda^2(1-\varphi_1^2)} + \frac{k_t}{\tau_t^2}}$$

$$\mu_{\vartheta_t(\text{post})} = \sigma_{\vartheta_t(\text{post})}^2 \left( \frac{\sum_{j=1}^{k_t} \mu_{jt}^*}{\tau_t^2} + \frac{(1-\varphi_1)\varphi_0 + \varphi_1 \vartheta_{t-1}}{\lambda^2(1-\varphi_1^2)} \right)$$

for  $1 < t < T$ :  $f(\vartheta_t | -) \propto \text{kernel of a } \mathcal{N}(\mu_{\vartheta_t(\text{post})}, \sigma_{\vartheta_t(\text{post})}^2)$  with

$$\sigma_{\vartheta_t(\text{post})}^2 = \frac{1}{\frac{1+\varphi_1^2}{\lambda^2(1-\varphi_1^2)} + \frac{k_t}{\tau_t^2}}$$

$$\mu_{\vartheta_t(\text{post})} = \sigma_{\vartheta_t(\text{post})}^2 \left( \frac{\sum_{j=1}^{k_t} \mu_{jt}^*}{\tau_t^2} + \frac{\varphi_1(\vartheta_{t-1} + \vartheta_{t+1}) + \varphi_0(1-\varphi_1)^2}{\lambda^2(1-\varphi_1^2)} \right)$$

for  $t = 1$ :  $f(\vartheta_t | -) \propto \text{kernel of a } \mathcal{N}(\mu_{\vartheta_t(\text{post})}, \sigma_{\vartheta_t(\text{post})}^2)$  with

$$\sigma_{\vartheta_t(\text{post})}^2 = \frac{1}{\frac{1}{\lambda^2} + \frac{\varphi_1^2}{\lambda^2(1-\varphi_1^2)} + \frac{k_t}{\tau_t^2}}$$

$$\mu_{\vartheta_t(\text{post})} = \sigma_{\vartheta_t(\text{post})}^2 \left( \frac{\varphi_0}{\lambda^2} + \frac{\varphi_1(\vartheta_{t+1} - (1-\varphi_1)\varphi_0)}{\lambda^2(1-\varphi_1^2)} + \frac{\sum_{j=1}^{k_t} \mu_{jt}^*}{\tau_t^2} \right) \quad (2.10)$$

- update  $\varphi_0$

$$f(\varphi_0 | -) \propto \text{kernel of a } \mathcal{N}(\mu_{\varphi_0(\text{post})}, \sigma_{\varphi_0(\text{post})}^2) \text{ with}$$

$$\sigma_{\varphi_0(\text{post})}^2 = \frac{1}{\frac{1}{s_0^2} + \frac{1}{\lambda^2} + \frac{(T-1)(1-\varphi_1)^2}{\lambda^2(1-\varphi_1^2)}}$$

$$\mu_{\varphi_0(\text{post})} = \sigma_{\varphi_0(\text{post})}^2 \left( \frac{m_0}{s_0^2} + \frac{\vartheta_1}{\lambda^2} + \frac{1-\varphi_1}{\lambda^2(1-\varphi_1^2)} \sum_{t=2}^T (\vartheta_t - \varphi_1 \vartheta_{t-1}) \right) \quad (2.11)$$

- update  $\lambda^2$

$$f(\lambda^2 | -) \propto \text{kernel of a } \text{invGamma}(a_{\lambda(\text{post})}, b_{\lambda(\text{post})}) \text{ with}$$

$$a_{\lambda(\text{post})} = \frac{T}{2} + a_\lambda$$

$$b_{\lambda(\text{post})} = \frac{(\vartheta_1 - \varphi_0)^2}{2} + \sum_{t=2}^T \frac{(\vartheta_t - (1-\varphi_1)\varphi_0 - \varphi_1 \vartheta_{t-1})^2}{2} + b_\lambda \quad (2.12)$$

- update  $\alpha$

if global  $\alpha$ :  $f(\alpha | -) \propto \text{kernel of a } \text{Beta}(a_{\alpha(\text{post})}, b_{\alpha(\text{post})})$  with

$$a_{\alpha(\text{post})} = a_\alpha + \sum_{i=1}^n \sum_{t=1}^T \gamma_{it} \quad b_{\alpha(\text{post})} = b_\alpha + nT - \sum_{i=1}^n \sum_{t=1}^T \gamma_{it}$$

if time specific  $\alpha$ :  $f(\alpha_t | -) \propto \text{kernel of a } \text{Beta}(a_{\alpha(\text{post})}, b_{\alpha(\text{post})})$  with

$$a_{\alpha(\text{post})} = a_\alpha + \sum_{i=1}^n \gamma_{it} \quad b_{\alpha(\text{post})} = b_\alpha + n - \sum_{i=1}^n \gamma_{it}$$

if unit specific  $\alpha$ :  $f(\alpha_i | -) \propto \text{kernel of a } \text{Beta}(a_{\alpha(\text{post})}, b_{\alpha(\text{post})})$  with

$$a_{\alpha(\text{post})} = a_{\alpha i} + \sum_{t=1}^T \gamma_{it} \quad b_{\alpha(\text{post})} = b_{\alpha i} + T - \sum_{t=1}^T \gamma_{it}$$

if time and unit specific  $\alpha$ :  $f(\alpha_{it} | -) \propto \text{kernel of a } \text{Beta}(a_{\alpha(\text{post})}, b_{\alpha(\text{post})})$  with

$$a_{\alpha(\text{post})} = a_{\alpha i} + \gamma_{it} \quad b_{\alpha(\text{post})} = b_{\alpha i} + 1 - \gamma_{it} \quad (2.13)$$

- update a missing observation  $Y_{it}$

*[Scrivere che è un altro step rispetto al MCMC  
di Pepe et al.]*

for  $t = 1$ :  $f(Y_{it} | -) \propto \text{kernel of a } \mathcal{N}(\mu_{Y_{it}(\text{post})}, \sigma_{Y_{it}(\text{post})}^2)$  with

$$\sigma_{Y_{it}(\text{post})}^2 = \frac{1}{\frac{1}{\sigma_{c_{it} t}^{2*}} + \frac{\eta_{1i}^2}{2\sigma_{c_{it+1} t+1}^{2*}(1-\eta_{1i}^2)}}$$

$$\mu_{Y_{it}(\text{post})} = \sigma_{Y_{it}(\text{post})}^2 \left( \frac{\mu_{c_{it} t}^* + \mathbf{x}_{it}^T \boldsymbol{\beta}_t}{\sigma_{c_{it} t}^{2*}} + \frac{\eta_{1i}(Y_{it+1} - \mu_{c_{it+1} t+1}^* - \mathbf{x}_{it+1}^T \boldsymbol{\beta}_{t+1})}{\sigma_{c_{it+1} t+1}^{2*}(1-\eta_{1i}^2)} \right)$$

for  $1 < t < T$ :  $f(Y_{it} | -) \propto \text{kernel of a } \mathcal{N}(\mu_{Y_{it}(\text{post})}, \sigma_{Y_{it}(\text{post})}^2)$  with

$$\sigma_{Y_{it}(\text{post})}^2 = \frac{1 - \eta_{1i}^2}{\frac{1}{\sigma_{c_{it} t}^{2*}} + \frac{\eta_{1i}^2}{\sigma_{c_{it+1} t+1}^{2*}}}$$

*Esempio: nuberi massimi e minimi?*

$$\mu_{Y_{it}(\text{post})} = \sigma_{Y_{it}(\text{post})}^2 \left( \frac{\mu_{c_{it}}^* + \eta_{1i} Y_{it-1} + \mathbf{x}_{it}^T \boldsymbol{\beta}_t}{\sigma_{c_{it}}^{2*}(1 - \eta_{1i}^2)} + \frac{\eta_{1i}(Y_{it+1} - \mu_{c_{it+1}t+1}^* - \mathbf{x}_{it+1}^T \boldsymbol{\beta}_{t+1})}{\sigma_{c_{it+1}t+1}^{2*}(1 - \eta_{1i}^2)} \right)$$

for  $t = T$ :  $f(Y_{it} | -)$  is just the likelihood of  $Y_{it}$  (2.14)

Finally, briefly highlight in Algorithm 1 the steps which compose the MCMC sampling algorithm. Regarding the computation of the fitting metrics LPML and WAIC, they follow classic ideas from [Chr+10] and [GHV13] respectively. ?? singler!

The core of the clustering process happens in the updating steps of  $\gamma_{it}$  and  $\rho_t$ . Their update step is indeed quite complex, and as we said before involves the check of compatibility issues. In any case, the general idea is that, for each unit  $i$  currently belonging to cluster  $j$ , we simulate to assign her to one of the existing clusters, plus to a new singleton cluster, and compute for each case the likelihood of this to happen, deriving probability weights to finally sample the decision for the next iteration. The key elements participating into the definition of such weights are the spatial cohesions and, with the JDRPM update, also the covariate similarities, which we will now both investigate.

me multiplica per le pmi!

## 2.2 Spatial cohesions analysis

The clustering procedure revolves around the product partition model (PPM). The simplest idea is to set  $P(\rho_t) \propto \prod_{j=1}^{k_t} C(S_{jt})$ , with the function  $C(S_{jt})$  that measures how tightly grouped the elements in  $A$  are considered to be. Then, to include spatial information, the idea is to extend the PPM from being a function of just  $C(S_{jt})$  to the more informed one  $C(S_{jt}, \mathbf{s}_{jt}^*)$ , where  $S_{jt}$  is the  $j$ -th cluster at time instant  $t$  and  $\mathbf{s}_{jt}^*$  is the subset of spatial coordinates of the units inside  $S_{jt}$ . For the sake of clarity, in this section where we are just interested in analysing the cohesions we employ the  $S_h$  notation to indicate a general  $h$ -th cluster, rather than the more pedantic  $S_{jt}$ .

Regarding the computation of spatial cohesion, several choices are available [PQ15]. The main common idea of the following formulas is to favour few spatially connected clusters rather than a lot of singleton ones, to derive more interpretable and meaningful results. For this reason, most of the cohesions employ the  $M \cdot \Gamma(|S_h|)$  term, which resembles the DP partitioning method that helps in reaching that goal.

We will now describe briefly all the cohesions which are implemented in the JDRPM model, and were implemented as well in the CDRPM model, and conduct tests on each of them, to see how the tuning of their parameters reflects on the computed values.

All the tests of Figures 2.3 and 2.4 refer to the partition of Figure 2.2, taken as test case here from a general fit on the same spatio-temporal dataset of Chapter 4. The following results, as well as the ones of the next section, are presented with the logarithm applied to better highlight the differences among them, otherwise for example all values could be really close together and make the analysis less understandable, and moreover because this is the actual perspective in which the