

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

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. ??, inglese!

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.

## 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 pedante  $S_{jt}$ . ??

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.

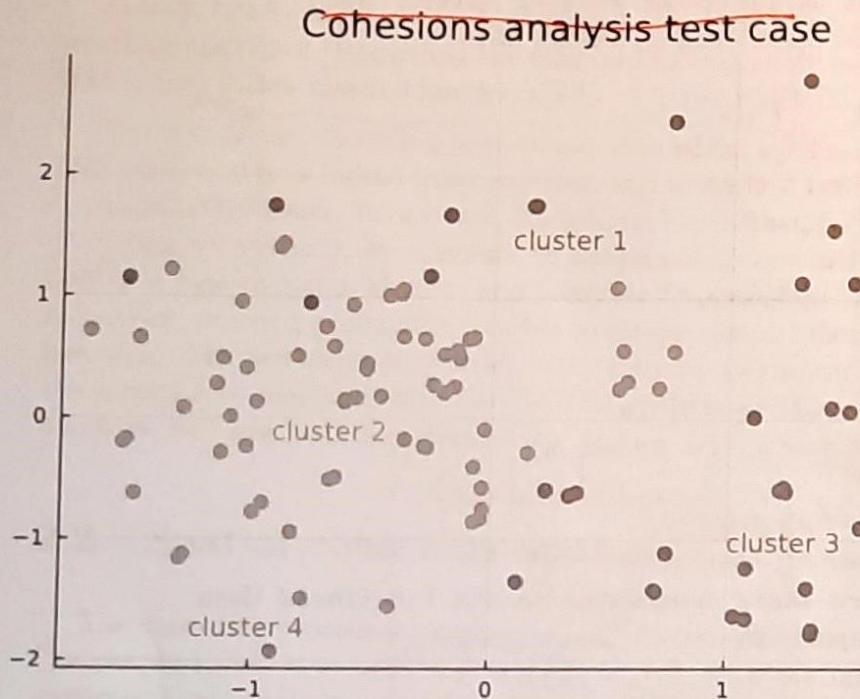
We will now describe briefly all the cohesions which are implemented in the JDRPM model, and were implemented as well in the CDRPM model, and conduct tests on each of them, to see how the tuning of their parameters reflects on the computed values. ?? quali? che valori? influence the

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

~ vuol dire che l'inglese è brutto

Prima spiega / introduce le coesioni (2.15), E..), poi spiega come sono nelle Figure 2.3 e 2.4.

implementation works. In fact, the fitting algorithm firstly saves the log-transformed values generated by the cohesions, in order to avoid numerical problems and instabilities, and secondly exponentiates and normalizes them into proper scaled probabilities, from which finally draw the cluster assignments.



? me  
numerico  
della qui?  
Questo sembra un "tutto"  
numerico.

Prima definizione  
di  $C_1$ ,  
per spiegare / commentare

Figure 2.2: Partition considered to analyse the spatial cohesions.

The first cohesion uses a tessellation idea from [DH01] that considers  $D_h = \sum_{i \in S_h} \|s_i - \bar{s}_h\|$  as the total distance from the units to the cluster centroid  $\bar{s}_h$ . The computation is then an adjustment of a decreasing function in terms of  $D_h$ , to give an higher weight on clusters which are denser, i.e. that have lower  $D_h$ , with an additional parameter  $\alpha$  to provide more control on the penalization.

$$C_1(S_h, s_h^*) = \begin{cases} \frac{M \cdot \Gamma(|S_h|)}{\Gamma(\alpha D_h) \mathbb{1}_{[D_h \geq 1]} + D_h \mathbb{1}_{[D_h < 1]}} & \text{if } |S_h| > 1 \\ M & \text{if } |S_h| = 1 \end{cases} \quad (2.15)$$

The second function provides, instead, a hard cluster boundary, where the weight is set to 1, i.e. 0 with the logarithm view, only if all the distances between all possible pairs of points inside the cluster are below the threshold parameter, i.e. if all units are "close enough" to each other. If this does not happen, even for a single pair of points, the returned value is 0, which corresponds to the maximum penalization since it would be  $-\infty$  in the logarithm perspective. The strictness of this requirement can be adjusted through the parameter  $a$ .

$$C_2(S_h, s_h^*) = M \cdot \Gamma(|S_h|) \cdot \prod_{i,j \in S_h} \mathbb{1}_{[\|s_i - s_j\| \leq a]} \quad (2.16)$$

Prima definizione di  $C_2$ , per i commenti / spiegazioni.

cosa che in queste figure? NON a' capisce!  
così è, il valore della funzione  $C_2$ .  
13

According to this function, from Figure 2.3 we can see how the purple cluster is considered the one with the highest cohesion, being a singleton. The runner-up is the green cluster, because it's the first among all the non-singletons which activates cohesion 2 when we increase the parameter  $a$ . The orange and blue clusters, instead, appear to be less dense since they require an higher value of  $a$  to "pass" the distance check.

However, cohensions  $C_1$  and  $C_2$  do not preserve the exchangeability property, meaning that if we would marginalize the random partition model over the last of  $m$  units we would not get to the same model as if we only had  $m - 1$  units. This coherence property, known as sample size consistency or addition rule [De +15], is instead often desirable, for theoretical or computational purposes, and the following two cohensions are able to provide it [MQR11].

Cohesion 3, called auxiliary similarity, treats the spatial coordinates  $s^*$  as if they were random, applying on them a model such as the Normal/Normal-Inverse-Wishart with  $\xi = (\mathbf{m}, V)$ ,  $s|\xi \sim N(\mathbf{m}, V)$  and  $\xi \sim NIW(\mu_0, \kappa_0, \nu_0, \Lambda_0)$ . The idea is to assign a larger weight on clusters which produce large marginal likelihood values, i.e. which according to the modelling are more probable to appear.

prima le def, el hin meni cose cosa oggi e so??

$$C_3(S_h, s_h^*) = M \cdot \Gamma(|S_h|) \cdot \int \prod_{i \in S_h} q(s_i | \xi_h) q(\xi_h) d\xi_h \quad (2.17)$$

??

On the same line there is cohesion 4, the double dipper cohesion [QMP15], which now employs the posterior predictive distribution rather than the prior predictive distribution of cohesion  $C_3$ .

$$C_4(S_h, s_h^*) = M \cdot \Gamma(|S_h|) \cdot \int \prod_{i \in S_h} q(s_i | \xi_h) q(\xi_h | s_h^*) d\xi_h \quad (2.18)$$

Another final idea comes from the cluster variance/entropy similarity function, a very general methodology to measure the closeness of a set of values which in fact will be used also for the covariates case. As in cohesion  $C_1$ , the idea is to derive a summary of the closeness of the units, summing the distance of the units from the cluster centroid, and then adjust the penalization through the parameter  $\varphi$ , to control how much penalize dissimilar values.

c'è lavoro?

$$C_5(S_h, s_h^*) = \exp \left\{ -\varphi \sum_{i \in S_h} \|s_i - \bar{s}_h\| \right\} \quad (2.19)$$

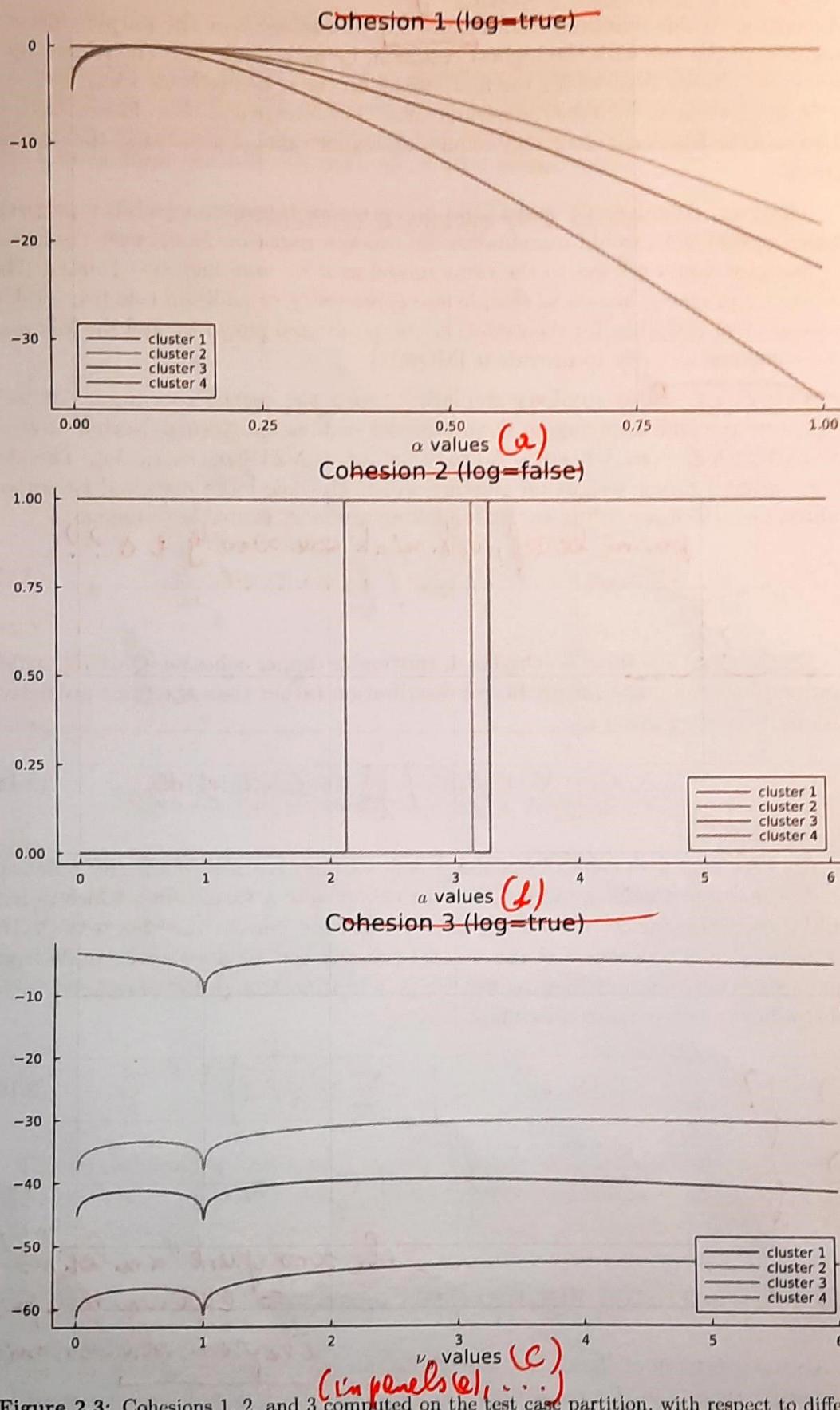
leჩե ենք?

$$C_6(S_h, s_h^*) = \exp \left\{ -\varphi \log \left( \sum_{i \in S_h} \|s_i - \bar{s}_h\| \right) \right\} \quad (2.20)$$

## 2.3 Covariates similarities analysis

A wide spectrum of choice is also available for covariates similarities [PQ18]. To account for them, the idea consists in extending the PPM to make it function of  $S_{jt}$ ,

le chiamiamo SIMILARITIES  
ma le forme toller le (2.21) le  
dove mettere ell'immagine del capitolo 2



**Figure 2.3:** Cohesions 1, 2, and 3 computed on the test case partition, with respect to different values of their tuning parameter. Cohesion 2 is without the logarithm applied just for plotting purposes, since otherwise the values would have been  $-\infty$  and 0.

$s_{jt}^*$ , and now also of  $X_{jt}^*$ , being this the  $p \times |S_{jt}|$  matrix storing the covariates of the units belonging to cluster  $j$  at time  $t$ , i.e.  $X_{jt}^* = \{\mathbf{x}_{it}^* = (x_{it1}, \dots, x_{itp})^T : i \in S_{jt}\}$ .

For the current implementation of JDRPM we decided to treat each covariate individually, therefore the PPM will actually be function of the set of unidimensional vectors which record the different  $p$  covariates of the units inside cluster  $S_{jt}$ , meaning  $\mathbf{x}_{jt1}^*, \dots, \mathbf{x}_{jtp}^*$ , of which all contributions will be considered independently one to each other. In this way, the final PPM form will be

$$P(\rho) \propto \prod_{h=1}^{k_t} C(S_{jt}, s_{jt}^*) \left( \prod_{r=1}^p g(S_{jt}, \mathbf{x}_{jtr}^*) \right) \quad (2.21)$$

where  $\mathbf{x}_{jtr}^*$  represents the vector recording the  $r$ -th covariate values for all the units inside cluster  $S_{jt}$ , i.e. row  $r$  of matrix  $X_{jt}^*$ . This choice is lighter from the theoretical and computational perspectives, and allows a mixture of numerical and categorical covariates without any problem. A unified and multidimensional consideration of the covariates is nonetheless possible, with proper adjustments on the similarities functions, and would have lead to a PPM of the form

$$P(\rho) \propto \prod_{h=1}^{k_t} C(S_{jt}, s_{jt}^*) g(S_{jt}, X_{jt}^*) \quad (2.22)$$

evidence i: "title" delle figure  
me scriveri ception complete

### Similarities analysis test case

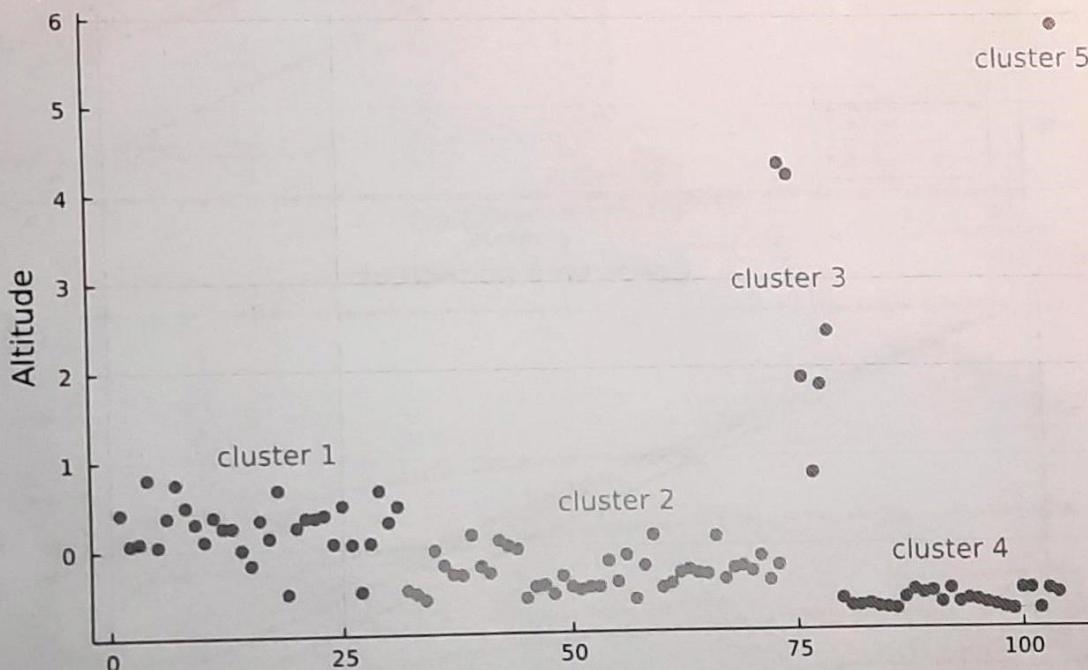


Figure 2.5: Partition considered to analyse the covariate similarities. The order of the units has been modified to plot together all units belonging to the same cluster.

As in the previous section, we will now discuss all the similarities implemented in the JDRPM model and perform tests on each of them. All the tests will be

In une lesie di skok shice, "test" vuol dire  
un test skok shice

?? copie questi  
test ??

referred to the test case partition displayed in Figure 2.5. Regarding the notation, we will again employ the simpler and general one dropping the spatio-temporal context.

The first similarity is the cluster variance/entropy similarity function, which as the name suggest can work with both numerical and categorical variables. The general form is

$$g_1(S_h, \mathbf{x}_h^*) = \exp \{-\varphi H(S_h, \mathbf{x}_h^*)\} \quad (2.23)$$

with  $H(S_h, \mathbf{x}_h^*) = \sum_{i \in S_h} (x_i - \bar{x}_h)^2$  for numerical covariates, being  $\bar{x}_h$  the mean value of the  $\mathbf{x}_h^*$  vector, and  $H(S_h, \mathbf{x}_h^*) = -\sum_{c=1}^C \hat{p}_c \log(\hat{p}_c)$  for categorical covariates, with  $\hat{p}_c$  indicating the relative frequency at which each factor appears. The parameter  $\varphi$  controls the amount of penalization to apply. This function can be extended quite easily to the multidimensional case, at least for the case of numerical covariates only, where the  $H$  function becomes  $H(S_h, X_h^*) = \sum_{r=1}^p \|\mathbf{x}_r - \bar{\mathbf{x}}_h\|$ .

Another popular choice is the Gower similarity function [Gow71], which was originally designed for a multivariate context, where vectors of covariates are compared to each other. As a consequence, some work was required to adapt it to the univariate case. The simple idea of comparing all cluster-specific pair-wise similarities leads to the total Gower similarity.

$$g_2(S_h, \mathbf{x}_h^*) = \exp \left\{ -\alpha \sum_{i,j \in S_h, i \neq j} d(x_i, x_j) \right\} \quad (2.24)$$

This function, however, is strictly increasing with respect to the cluster size, meaning that it will naturally tend to propose a large number of small clusters. For that reason, a correction can be applied, accounting for the size of the cluster  $S_h$ , leading to the average Gower similarity.

$$g_3(S_h, \mathbf{x}_h^*) = \exp \left\{ -\frac{2\alpha}{|S_h|(|S_h| - 1)} \sum_{i,j \in S_h, i \neq j} d(x_i, x_j) \right\} \quad (2.25)$$

In both functions,  $d(x_i, x_j)$  is the Gower dissimilarity between  $x_i$  and  $x_j$ , with  $d(x_i, x_j) = |x_i - x_j|/R$  in the case of numerical covariates, being  $R = \max(\mathbf{x}) - \min(\mathbf{x})$  the range of the covariate values, considering all the units independently from their cluster, while  $d(x_i, x_j) = \mathbb{1}_{[x_i \neq x_j]}$  in the case of categorical covariates. This is a dissimilarity since values closer to 0 refer to similar data, while values closer to 1 to dissimilar data; therefore the minus sign inside  $g_2$  and  $g_3$  exponents converts the function to a similarity. For the multidimensional design there are natural extensions of the  $d(\mathbf{x}_i, \mathbf{x}_j)$  function.

The last similarity recalls the structure of the spatial cohesion  $C_3$ , where now are the covariates to be treated as if they were random variables. However, since we chose to deal with each covariate individually, in this unidimensional setting a Normal/Normal-Inverse-Gamma model is employed, with  $\boldsymbol{\xi} = (\mu, \sigma^2)$ ,  $x|\boldsymbol{\xi} \sim \mathcal{N}(\mu_0, \sigma^2)$ , and  $\mu \sim \mathcal{N}(\mu_0, \sigma^2/\lambda_0)$ ,  $\sigma^2 \sim \text{invGamma}(a_0, b_0)$ , i.e.  $\boldsymbol{\xi} \sim \mathcal{N}\text{invGamma}(\mu_0, \lambda_0, a_0, b_0)$ . A multivariate extension is thus possible through the same modelling style of the spatial coordinates case.

$$g_4(S_h, \mathbf{x}_h^*) = \int \prod_{i \in S_h} q(x_i | \boldsymbol{\xi}_h) q(\boldsymbol{\xi}_h) d\boldsymbol{\xi}_h \quad (2.26)$$

*Here*

All the similarities, except for  $g_2$ , agree to classify the non-singleton clusters with the red being the one with the highest similarity, followed by the orange, blue, and green, as we can see from Figures 2.6 and 2.7. Intuitively, Figure 2.5 seems to confirm this ranking, by visually reasoning about the sparsity of each cluster. Similarities  $g_4$  and  $g_2$  seem to be the only ones which are able to give a considerable weight also to the green partition, which is indeed sparser but collects all the large, sort of outliers, values of the covariate at test.

Regarding their flexibility, we tested their behaviour changing the testing covariate. Similarity  $g_1$  appeared to be the most adaptive one, working well (i.e. returning very reasonable orders in the clusters). However, it tends to penalize a lot sparse clusters, as we saw in the previous test case. This could suggest to maybe implement other distance metrics rather than the  $L^2$  norm in the computation of  $H(S_h, \mathbf{x}_h^*)$ . On the other hand, similarity 4 appeared to be quite sensible to the parameters regulating the invGamma distribution for  $\sigma^2$ , suggesting that each covariate that we want to include in the clustering process should have her properly-tuned pair of  $a_0$  and  $b_0$  parameters. The JDRPM implementation will provide this flexibility in assigning separately those parameters for each clustering covariate.

{  
single  
hypo  
flexible}

Prime di parlare delle ottimizzazioni del codice, dove  
 brevemente spieghi che il calcolo delle posteriori per i modelli  
 bayesiani si fa con metodi di simulazione, gli MCMC. Spiega in  
 3 fasi cosa sono, e descrive brevemente quello di Pepe et al.  
 Chapter 3 → queste parole le potrebbe mettere nel Cap precedente, prima  
 di fare descrivere gli algoritmi. In questo,  
 qui si fa riferimento a quelle parole

## Implementation and optimizations

D Gliogliamo l'algoritmo MCMC per calcolare la posteriori del nostro  
 modello e stiamo scrivendo in Julia...

To implement our updated model, and to do it efficiently, we decided to opt for the Julia language [Bez+17].

Julia is a relatively new programming language that combines the ease and expressiveness of high-level languages to the efficiency and performance of low-level ones. Such balance largely achieved through the just-in-time (JIT) compilation, implemented using the LLVM framework, together with many nice features like dynamic multiple dispatch and heavy code specialization against multiple types. This design choice allows Julia to be used interactively, in the same fashion to the R, Matlab, or Python consoles, while also supporting the more traditional program execution style of statically compiled languages such as C, C++ and Fortran. This solution provides faster development phases, since for examples code sections can be evaluated and tested line by line, guaranteeing and at the same time efficient implementations. Performances are further enhanced by the native integration of the optimized BLAS [Law+79] and LAPACK libraries for linear algebra operations, which are often at the core of all scientific applications. Regarding productivity, instead, Julia provides an extensive ecosystem, currently comprised by more than ten thousand packages, spanning over almost all branches of science and engineering. Moreover, most of them are already highly tested and optimized, and can therefore reduce the implementation time required to the users. As an example, in this work we employed the `Distributions` [Bes+21] [Lin+19] and `Statistics` packages, while the original C implementation had to write all the statistical functionalities from scratch. Considering all this, the Julia choice was deemed very natural.

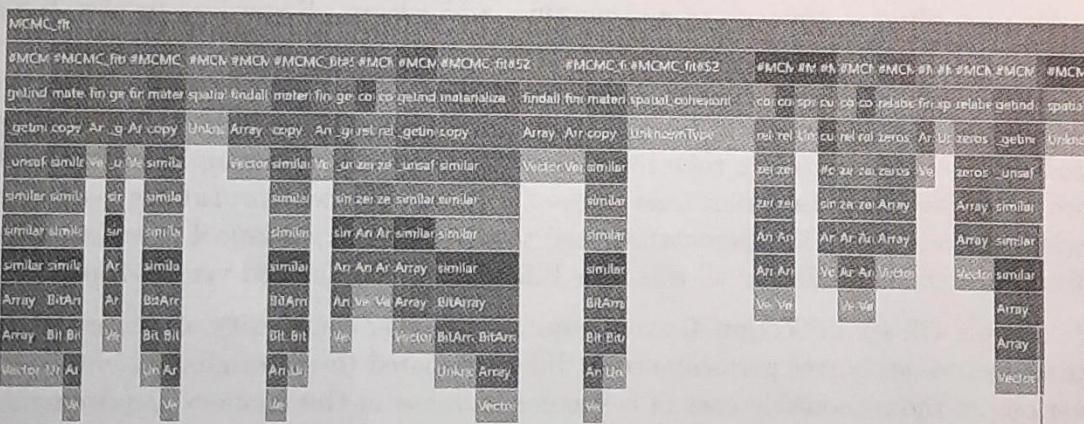
As we will see in Section 4, even despite the higher complexity of the model we managed to get better performance in Julia compared to the original C implementation, at the reasonable cost of a smaller increase in the memory requirements, which nowadays should not be a big deal. Anyway, this improvement was possible through several refinements and optimizations, which we will now discuss.

### 3.1 Optimizations

One of the major issues encountered during the design of the fitting function in Julia was controlling the amount of memory and allocations that some functions, structures, or algorithms would require. At the beginning of the development and testing, where the correctness of the algorithm was the only priority, we saw in fact that most of the execution time was actually spent by the garbage collector of Julia, which had the burden of track all the allocated memory and reclaim the unused one in order to make it available again for new computations.

Together with avoiding useless allocations, or in general managing more precisely the memory, another important aspect to focus on to get better performance is ensuring type stability of the function. Luckily, Julia disposes of several tools to inspect and address both problems.

Regarding the second point, there are several packages such as Cthulhu or even the simple `@code_warntype` macro which allow to ensure that a function is type-stable, i.e. that the types of all variables can be correctly predicted by the compiler and stay consistent throughout the execution. Type stability is crucial for performance as it enables the compiler to generate optimized machine code, removing the load derived from dynamic type checks. In fact, Julia is dynamically typed, which means that variables do not have to be necessarily declared with their type, in contrast to fully statically typed languages such as C, and they could change it during execution. For example, a variable initialized as an integer could later become a float, or even a string. However, for performance purposes, these dynamicisms should be avoided, and those tools help to ensure that this happens. We managed to reduce all the useless type instabilities, but some are instead still present to guarantee some flexibility to the users, e.g. to select at run time the cohesion and similarity functions to use in the fit.



**Figure 3.1:** Flame graph derived from a simple test fit with  $n = 20$ ,  $T = 50$ , ran for 10000 iterates. Each section, identified with a colored box, represents a function call with its stack trace below, i.e. all the other subsequent function calls generated from the initial top one. The widths are proportional to allocation counts, i.e. to how many times a memory allocation was required for their corresponding box. The plot should be read from top to bottom with respect to function calls, and from left to right with respect to the fitting steps.

questo per le varie misure nel testo

```

aux1 = k0 * sdim; aux2 = k0 + sdim
Psi_n = Psi .+ S .+ aux1 / (aux2) .* Mtmp

```

This final version exploits the `StaticArrays` package of Julia, which allows to use vectors and matrices more efficiently if their size is known at compile time; which is the case of the spatial cohesions computation since working with planar spatial coordinates we will always have  $2 \times 1$  vectors and  $2 \times 2$  matrices. The benefits of this final version are that we maintain the natural mathematical form of the first one, improving the clearness of the code, together with the efficiency of the second one, since now with static structures the compiler is able to optimize all memory allocations as it was doing with the simple scalar variables.

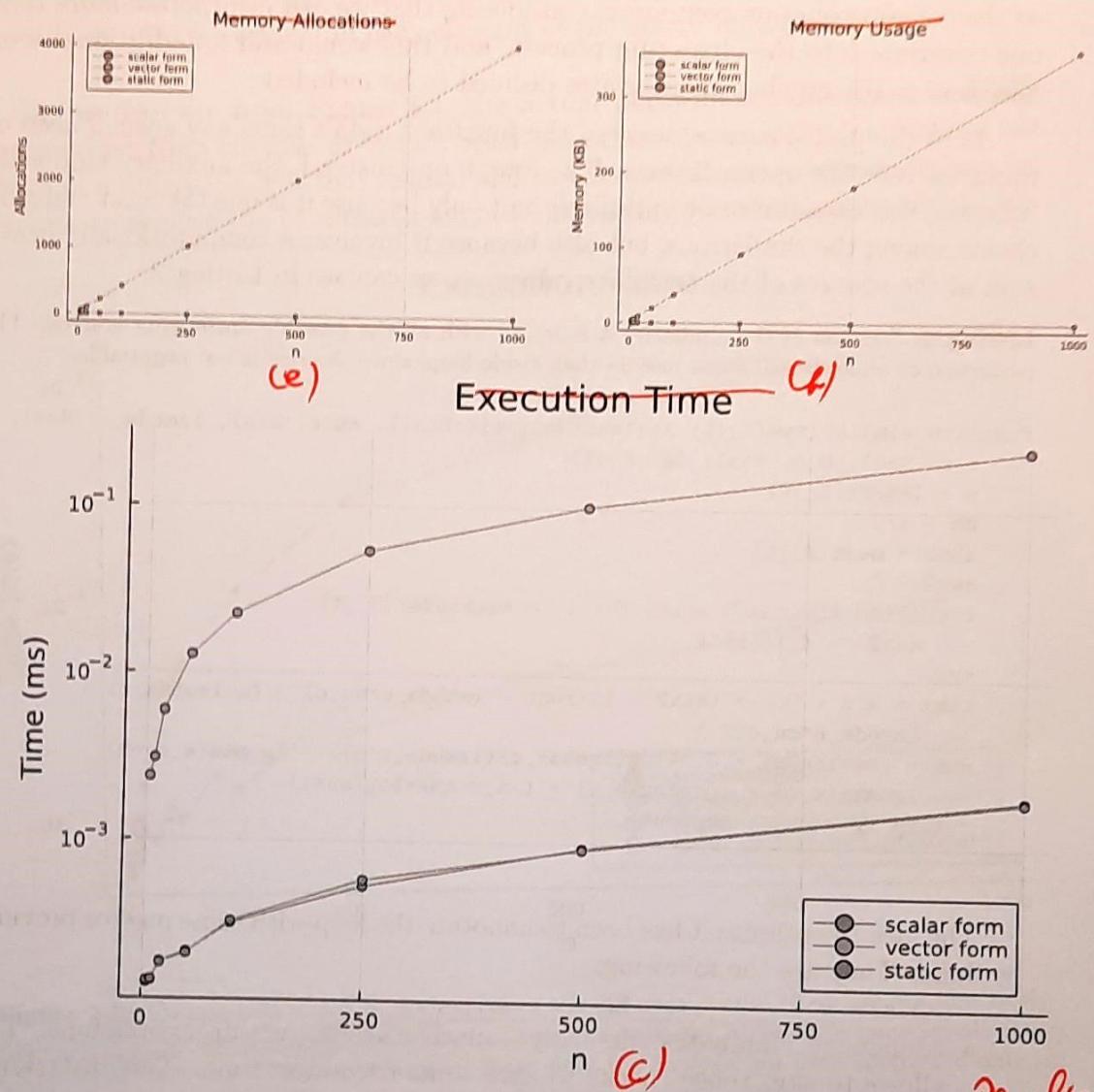


Figure 3.2: Performance comparison between the three versions of the cohesion 4 function. Tests ran through the `BenchmarkTools` package of Julia by randomly generating the spatial coordinates of the various test sizes  $n$ , with similar results standing for cohesion 3. The memory allocation and usage plots are constant at zero for both the scalar and vector static cases.

Penalized for  
Handwritten notes  
??

Figure 3.2 shows the comparison of their performances, where we can see how