

DRAFT

LAPO SANTI

CONTENTS

1. INTRODUCTION

In the present work, we specify and estimate a model for pairwise or binary data capable of clustering the data within blocks. The novelty lies in the fact that these blocks are ordered, meaning that they present a hierarchical relationship of some sort.

We start by presenting the unordered model, which is essentially a traditional Stochastic Block Models (SBMs). Pairwise data are modelled as binomially distributed, where the probability of success depends both on the membership of the two pair members, stored within the parameter z , and on a blocks-related probabilities matrix P .

The classical SBM models P with a Beta(1,1) prior. In the present work instead we model this P matrix in different ways in order to induce an ordering among the clusters.

We start by reviewing the unordered model, then, by drawing inspiration from the literature on image recognition, we start representing the matrix P as an ordered poset, meaning a discrete set where we impose an ordering structure. This structure is strongly influenced by the properties that we would like our ranked cluster to have.

For example, if we work under the Strong Stochastic Transitivity (SST) axiom, we are asking that if the individual i is ranked higher than j and j is ranked higher than q , then i must be ranked higher than the higher between i and j . There are other axioms that we might choose, like the Weak Stochastic Transitivity one (WST), which is less stringent than the SST, or the Linear Stochastic Transitivity one (LST), which encompasses both the two aforementioned ones.

However, assuming the SST axiom and imposing it over the poset defined on P lets emerge an interesting structure, meaning that of the Level Sets $L_{(k)}$. Each Level Set is defined as the set of entries on a specific diagonal of the upper triangular matrix P , starting from $L_{(0)}$ which is the main diagonal, $L_{(1)}$ is the diagonal above and so on, until $L_{(K-1)}$. We impose that $L_{(0)}$ is always equal to $\frac{1}{2}$, since we assume that individuals within the same block have equal probability of being preferred.

We proceed to specify over each Level Set a prior probability distribution, which governs the probability law of the entries within that diagonal. The probability distribution is a truncated normal distribution $TRUNCNORM(\mu_{(k)}|\alpha, \beta_{\max}, \sigma)$.

Each distribution is endowed with a specific and increasing $\mu_{(k)}|\alpha, \beta_{\max}$, which effectively induces a hierarchical model, given that α, β_{\max} are parameters common to each level set. While α specifies the rate of increase of the Level Sets' means $\{\mu_{(k)}, k = 0, \dots, K-1\}$, β_{\max} specifies the maximum attainable probability within the matrix P .

The last and most important hyperparameter σ is common to each distribution over the level sets. As σ increases, the truncated normals become more and more flat, and therefore we are back to a case where the entries of P are not distributed in a way that is not significantly different from a uniform.

This model should therefore encompass the unordered one, for values of σ that are large enough.

2. LITERATURE REVIEW

The present work lies at the intersections of at least three well defined streams of literature. The first one is the prolific literature on Stochastic Block Models. The second one is the ranking literature, based on the Bradley Terry model. The final one is the partial order literature.

3. THE UNORDERED MODEL SPECIFICATION

This is a model for pairwise count data. We explicitly model the results of the interactions between two individuals i and j . Given N observations, the likelihood is

$$(1) \quad p(y|z, P, K) = \prod_{i=2}^{N-1} \prod_{j=i}^N p(y_{ij}|z, P, K)$$

$$(2) \quad = \prod_{i=2}^{N-1} \prod_{j=i}^N \binom{n_{ij}}{y_{ij}} p_{z_i, z_j}^{y_{ij}} (1 - p_{z_i, z_j})^{n_{ij} - y_{ij}}$$

where n_{ij} denotes the total number of interactions between the two individuals i and j and y_{ij} is the number of successes of the individual i in interacting with j . The probability of success is given by p_{z_i, z_j} which consists of two parameters. The $K \times K$ matrix P and the $N \times 1$ vector z .

The vector z has entries z_i taking values over the discrete and finite set $\{1, \dots, K\}$, and it is an indicator variable such that if $z_i = k$ individual i belongs to block k .

The matrix P contains the probabilities of success for individuals belonging to each possible blocks combination. For this reason P is $K \times K$. Therefore, the parameter p_{z_i, z_j} consists in the probability of success in an interaction between one individual belonging to block z_i and another of block z_j .

3.1. Prior Specification. Starting with the parameter P , we assume that its entries, namely $p_{k,k'}$, are independent and identically $Beta(a, b)$ distributed random variable. By setting $a = b = 1$ they collapse to a uniform distribution.

$$(3) \quad p_{k,k'} \sim Beta(1, 1) \quad \text{for } k, k' = 1, \dots, K$$

Second, we assume that the z_i 's are independent and identically drawn from a multinomial distribution with one trial and probability vector $(\theta_1, \dots, \theta_K)$. We can write then:

$$(4) \quad z_i | \boldsymbol{\theta} \sim \text{Multinomial}(1, \boldsymbol{\theta}) \quad \text{for } i = 1, \dots, N$$

To have more flexibility in the blocks sizes, we put an hyper-prior on the $\theta_1, \dots, \theta_K$, assuming that they are drawn from a Dirichlet distribution with parameter the $K \times 1$ vector $\boldsymbol{\gamma}$.

By marginalizing out θ , following the common practice in the literature, we can express the marginal distribution of z as:

$$(5) \quad p(\mathbf{z}|\boldsymbol{\gamma}) = \frac{\Gamma(\sum_{k=1}^K \gamma_k)}{\prod_{k=1}^K \Gamma(\gamma_k)} \frac{\prod_{k=1}^K \Gamma(n_k + \gamma_k)}{\Gamma(\sum_{k=1}^K (n_k + \gamma_k))}$$

where n_k is the number of players assigned to block k .

Finally, we assume that the number of clusters K follow a Poisson distribution $\text{Poisson}(\lambda = 1)$, subject to the condition $K > 0$.

4. CONNECTION WITH IMAGE RECOGNITION

By drawing inspiration from the literature in Image Recognition, in particular an article of Noel Cressie and Jennifer Davidson, we represent our matrix P as if it was an image, its probabilities as if they were pixels of different intensities, and its entries' indices as if they were pixels' locations.

Let us denote a generic entry of P as a vector s in \mathbb{N}^2 . The quantity $Z(s)$ denotes the probability value at the entry index s . We rewrite the whole matrix as

$$(6) \quad Z = \{Z(s) : s \in D\}$$

where D is the set of entries' indices of the matrix P , that is:

$$(7) \quad D = \{(u, v) : u = 1, \dots, K; v = 1, \dots, K\}.$$

Now, let us consider a temporal Markov process $\{Z(t) : t=1,2,\dots\}$. The Markov property can be generalised from a one-dimensional time-process to a two-dimensional space with both a conditional and a joint specification.

We draw a connection between the class of models called Partially Ordered Markov models which allows us to efficiently compute the joint probabilities of the prior on P , and the literature on preference learning.

Then we introduce the notion of partial order among the P entries. Let's take once more D . The binary relation \succeq on D is said to be a partial order if For any $x \in D$, $x \prec x$ (reflexivity). For any $x, y, z \in D$, $x \prec y$ and $y \prec z$ implies $x \prec z$ (transitivity). For any $x, y \in D$, $x \prec y$ and $y \prec x$ implies $x = y$ (antisymmetry). Then we call (D, \prec) a partially ordered set, or a poset. For example, the set of all subsets of a given set, with the relation \prec being set inclusion, is a poset.

Regarding the matrix P , we can check whether and how to use the definition of Poset under the three different axiomatic frameworks that specify different (stochastic) transitivity requirements.

- (1) Weak Stochastic Transitivity (WST): $\mathbb{P}(x \prec y) \geq \frac{1}{2}$ and $\mathbb{P}(y \prec z) \geq \frac{1}{2}$ imply $\mathbb{P}(x \prec z) \geq \frac{1}{2}$, for all $x, y, z \in \mathcal{A}$.
- (2) Strong Stochastic Transitivity (SST): $\mathbb{P}(x \prec y) \geq \frac{1}{2}$ and $\mathbb{P}(y \prec z) \geq \frac{1}{2}$ imply $\mathbb{P}(x \prec z) \geq \max\{\mathbb{P}(x \prec y), \mathbb{P}(y \prec z)\}$, for all $x, y, z \in \mathcal{A}$.
- (3) Linear Stochastic Transitivity (LST): $\mathbb{P}(a \prec b) = F(\mu(a) - \mu(b))$ for all $a, b \in \mathcal{A}$, where $F : \mathbb{R} \rightarrow [0, 1]$ is an increasing and symmetric function (referred to as a "comparison function"), and $\mu : \mathcal{A} \rightarrow \mathbb{R}$ is a mapping from the set \mathcal{A} of alternatives to the real line (referred to as a "merit function").

Each of these axioms, produces a different P structure. Assuming, without loss of generality, that block 1 is the strongest, and by imposing the main diagonal to be equal to $\frac{1}{2}$ we can visualise a matrix following WST as:

$$(8) \quad P^{WST} = \begin{pmatrix} p_{1,1} & p_{1,2} & \dots & p_{1,K} \\ p_{2,1} & p_{2,2} & \dots & p_{2,K} \\ \vdots & \vdots & \vdots & \vdots \\ p_{K,1} & p_{K,2} & \dots & p_{K,K} \end{pmatrix} = \begin{pmatrix} 1/2 \leq & p_{1,2} & \dots & p_{1,K} \\ p_{2,1} \leq & 1/2 \leq & \dots & p_{2,K} \\ \vdots & \vdots & \vdots & \vdots \\ p_{K,1} & p_{K,2} & \dots & 1/2 \end{pmatrix}$$

Instead, under SST, we would observe:

$$(9) \quad P^{SST} = \begin{pmatrix} p_{1,1} & p_{1,2} & \dots & p_{1,K} \\ p_{2,1} & p_{2,2} & \dots & p_{2,K} \\ \vdots & \vdots & \vdots & \vdots \\ p_{K,1} & p_{K,2} & \dots & p_{K,K} \end{pmatrix} = \begin{pmatrix} 1/2 \leq & p_{1,2} \leq & \dots & p_{1,K} \\ \vee | & \vee | & & \vee | \\ p_{2,1} \leq & 1/2 & \dots & \leq p_{2,K} \\ \vdots & \vdots & \vdots & \vdots \\ p_{K,1} \leq & p_{K,2} \leq & \dots & \leq 1/2 \end{pmatrix}$$

With regard of LST , it is a generalisation of the two axioms and therefore it includes both aforementioned cases (??) and (??), depending how one specifies F and μ . Given the LST definition, we can calculate p_{ij} as follows:

$$(10) \quad p_{ij} = F(\mu(i) - \mu(j))$$

where i and j range from 1 to K and represent the alternatives in the set \mathcal{A} .

All the three axiomatic frameworks satisfy (in the LST case we need to check the functional form of F and μ) of the three conditions for being a poset. And therefore, we take advantage of the poset structure.

We can now describe the correspondence referred to above. This connection opens up a large literature on graphical models, outside of statistical image analysis, that we return to in Section 6.2. Let (D, F) be a directed acyclic graph, where $D = \{y_1, \dots, y_n\}$, a finite set. To construct a poset to which this digraph corresponds, we define the binary relation \prec on D by

$$y_i \prec y_i, \text{ for } i = 1, \dots, n;$$

$$y_i \prec y_j, \text{ if there exists a directed path from } y_i \text{ to } y_j \text{ in } (D, F).$$

Notice that several different directed acyclic graphs can yield the same poset. Conversely, given a finite poset (D, \prec) , a corresponding directed acyclic graph can be obtained by defining the set of edges F as follows: $(y_i, y_j) \in F$ if and only if $y_i \prec y_j$ and there does not exist a third element

$$z \neq y_i, y_j \text{ such that } y_i \prec z \prec y_j.$$

We saw above that the correspondence is many-to-one. Given a finite poset, one may construct a class of directed acyclic graphs; the correspondence described above is in a sense the minimal directed acyclic graph since it has the smallest possible directed edge set. However, if one starts with a directed acyclic graph, the corresponding poset is unique.

From the point of view of image modeling, we are more interested in the directed-acyclic-graph description because we are able to specify directly the spatial relations between pixel locations.

Let us introduce the notion of level set (also known as indifference set [ref:shah2016]). Imagine to partition the $\frac{K*(K-1)}{2} + K$ elements of the upper triangular P matrix into the union of K disjoin level sets $\{L_{(k)}\}_{k=0}^{K-1}$ of sizes $|L_{(k)}|$ such that $\sum_{k=0}^{K-1} |L_{(k)}| = \frac{K*(K-1)}{2} + K$. We write that the pair $(i, j) \sim (i', j')$ if they belong to the same level set.

$L_{(0)}$ corresponds to the main diagonal, and it has size K . $L_{(1)}$ corresponds to the diagonal above the main one, and it has size $K - 1$, and so on up to $L_{(K-1)}$, which just a single element, corresponding to the upper-right entry of the matrix P .

We say that a matrix P' respects the level set partition $\{L_{(k)}\}_{k=0}^{K-1}$ if

$$(11) \quad p(P') = p(P) \text{ for all quadruples } (i, j, i', j') \text{ such that } i \sim i' \text{ and } i \sim j'$$

In the literature we typically have the level sets defined directly over the p_{ij} that is, the probability of individual i to be preferred over individual j , without any block partition. If instead, a block partition is introduced, this satisfies the definition of level set (??), and it has a very natural interpretation in the context of ranking. For instance, in buying cars, frugal customers may be indifferent between high-priced cars; or in ranking news items, people from a certain country may be indifferent to the domestic news from other countries

What we are doing here is somewhat a grouping of the blocks, which as we said can be seen as level sets, again into a higher-tier level sets. The interpretation is not as straightforward, but it induces a kind of regularity among the relations between blocks, meaning that the probability of block 1 to be preferred to block 2 is drawn from the same distribution of the probability of block 2 being preferred to block 3, since they belong to the same diagonal of P , i.e. the same level set.

5. THE WST MODEL

This model implements the Weak Stochastic Transitivity assumption.

To have the WST axiom satisfied, we need the upper triangular entries of the P matrix to be always greater than 0.5

Therefore, to implement this model a Bayesian setting, we must change little with respect to the Unordered model.

We leave everything the same, but we force the upper triangular entries to be greater than 0.5. We can enforce effectively this condition by modifying the proposal distribution, and also the prior distribution.

6. THE SST MODEL

We want to model our matrix P according to the SST axiom. If we assume that block 1 is the first ranked block, meaning the one having the highest chance of being preferred versus the others, this implies that the upper triangular entries of P must be increasing as we move from the main diagonal towards the top-right corner. At the same time, we want some flexibility, in order to allow for the some variability within the same level set.

In (??) we provide an example of a $P_{4 \times 4}^{SST}$ matrix, that is a connection probability matrix for $K = 4$ which satisfies the SST condition.

$$(12) \quad \left[\begin{array}{cccc|c} p_{11} & \leq & p_{1,2} & \leq & p_{1,3} & \leq & p_{1,4} \\ \text{V}\mid & & \text{V}\mid & & \text{V}\mid & & \text{V}\mid \\ 1 - p_{1,2} & \leq & p_{2,2} & \leq & p_{2,3} & \leq & p_{2,4} \\ \text{V}\mid & & \text{V}\mid & & \text{V}\mid & & \text{V}\mid \\ 1 - p_{1,3} & \leq & 1 - p_{2,3} & \leq & p_{3,3} & \leq & p_{3,4} \\ \text{V}\mid & & \text{V}\mid & & \text{V}\mid & & \text{V}\mid \\ 1 - p_{1,4} & \leq & 1 - p_{2,4} & \leq & 1 - p_{3,4} & \leq & p_{4,4} \end{array} \right] \quad \begin{array}{l} L_{(3)} := \{p_{ij} \mid j - i = 3\} \\ L_{(2)} := \{p_{ij} \mid j - i = 2\} \\ L_{(1)} := \{p_{ij} \mid j - i = 1\} \\ L_{(0)} := \{p_{ij} \mid j - i = 0\} \end{array}$$

The entries $p_{ij} \in L_{(k)} \subset P$, where $L_{(k)}$ denotes a level set, satisfy the condition $j - i = k$. In other words, they correspond to the entries on the k -th diagonal above the main diagonal in the upper triangular matrix P . Entries within the same level set $L_{(k)}$ are bounded both above and below by the entries in $L_{(k+1)}$ and $L_{(k-1)}$ respectively. Additionally, we set the probabilities within the main diagonal to be 0.5, reflecting an equal chance of preference if two items belong to the same block. The maximum attainable probability is determined by the user and denoted as β_{\max} . This constraint is imposed to prevent extreme probabilities, such as values near 0 or 1.

Lastly, we assume that $p_{ij} \in L_{(k)}$ for $k = 1, \dots, K - 1$ are identically and independently distributed according to a truncated normal distribution. The truncations are set at 0.5 and β_{\max} , with mean $\mu_{(k)}$ and variance σ^2 denoted as parameters:

$$(13) \quad p_{ij}^{(k)} \sim \text{TruncatedNormal}(\mu_{(k)}, \sigma^2; 0.5, \beta_{\max})$$

where $0.5 < \mu_{(1)} < \mu_{(2)} < \dots < \mu_{(K-1)} < \beta_{\max}$, so that the inequality constraints are respected. Going back to the $P_{4 \times 4}^{SST}$ matrix example, the prior structure will look as follows:

$$(14) \quad \left[\begin{array}{cccc|c} p_{11} & p_{12} & p_{13} & p_{14} & p_{ij}^{(3)} \sim \text{TruncatedNormal}(\mu_{(3)}, \sigma^2; 0.5, \beta_{\max}) \\ 1 - p_{12} & p_{22} & p_{23} & p_{24} & p_{ij}^{(2)} \sim \text{TruncatedNormal}(\mu_{(2)}, \sigma^2; 0.5, \beta_{\max}) \\ 1 - p_{13} & 1 - p_{23} & p_{33} & p_{34} & p_{ij}^{(1)} \sim \text{TruncatedNormal}(\mu_{(1)}, \sigma^2; 0.5, \beta_{\max}) \\ 1 - p_{14} & 1 - p_{24} & 1 - p_{34} & p_{44} & p_{ii}^{(0)} = 0.5 \end{array} \right]$$

To enforce the constraint $0.5 < \mu_{(1)} < \mu_{(2)} < \dots < \mu_{(K-1)} < \beta_{\max}$, we consider a power-law function $y_{(k)} = 0.5 + g(k)^\alpha \mathbb{I}(0.5, \beta_{\max})$. This power law is monotonically increasing in $g(k)$, and it is flexible enough to account for different increasing rates by modulating accordingly the α values.

- For α values exceeding 1, a convex power-law function emerges, engendering a steady but accelerating increase in the level set truncations toward β_{\max} .
- When α is between 0 and 1, the power-law function becomes concave, promptly pushing values toward β_{\max} . This reflects a pronounced bias toward higher probabilities.
- Notably, as α is equal to 1, the power-law function becomes linear, leading to a constant increment in the $\mu_{(k)}$ parameter.

We have that

$$(15) \quad \mu_{(k)} = \frac{y_{(k)} + y_{(k+1)}}{2} \quad \text{where}$$

$$(16) \quad y_{(k)} = \left(\frac{(\beta_{\max} - 0.5)^{(1/\alpha)}}{K} \times k \right)^{\alpha} + 0.5 \quad \text{for } k = 0, \dots, K$$

We place a uniform hyper-prior on the parameter α in order to estimate this important parameter, upon which the means of the level sets entries hinge:

$$\alpha \sim \text{Uniform}(0, 3)$$

We also place a uniform hyper-prior on the parameter σ^2 , which stays constant across the level sets:

$$\sigma^2 \sim \text{Uniform}(0, 1)$$

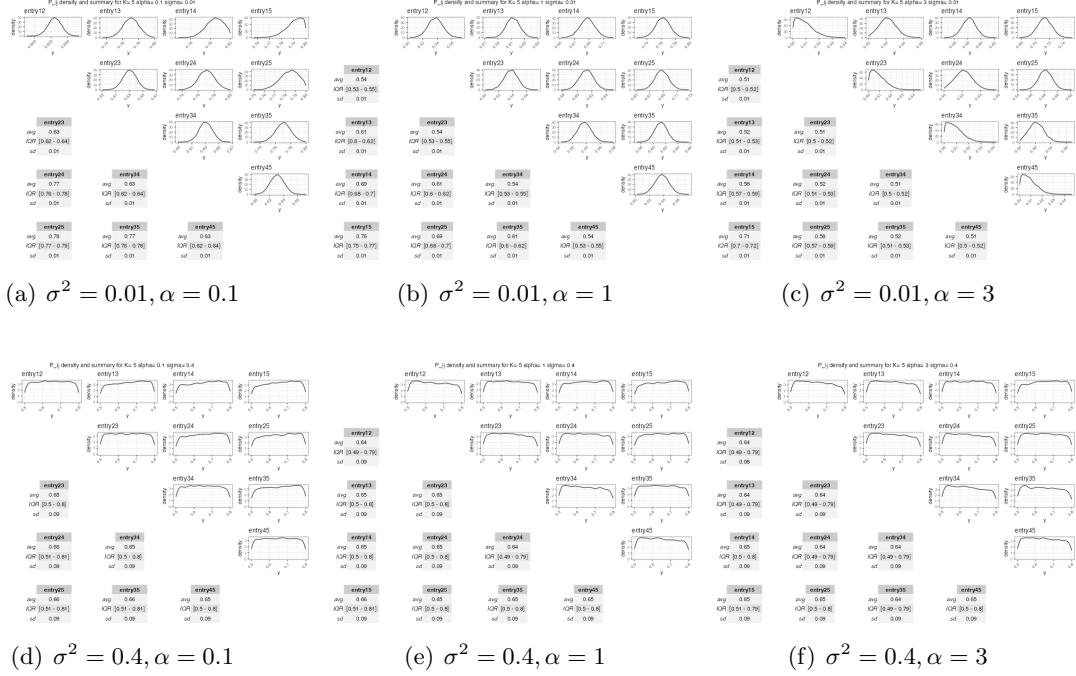


FIGURE 1. In the upper triangular side of each box above, we have the density of the p_{ij} , estimated over $n = 1000$ samples generated according to the combination of parameters σ^2, α reported below each box. K is set to 5. In the lower triangular part of each box, we have tables with some summary statics, namely the mean, the interquartile range, and the standard deviation.

6.1. Heteroskedastic variance. In the prior above, we are assuming that each entry of the P matrix has the same dispersion around the mean. However, in applications like Tennis, we witness a self-reinforcing mechanism according to which stronger players tend to play more. Therefore, we collect more data on them, and our estimates are more precise. This phenomenon does not need to be circumstanced just to Tennis. We may also think about certain goods being preferred by customers, located in the best position in the supermarket, and therefore being bought ever more. This would also lead to less data on the less preferred goods, and therefore, to higher variance.

To model also this aspect of the reality, we introduce another hyper parameter, namely ϕ , which allows the model to interpolate between a perfectly homoskedastic model to a perfectly heteroskedastic one.

The new distribution of the p_{ij} will look like:

$$(17) \quad p_{ij}^{(k)} \sim \text{TruncatedNormal}(\mu_{(k)}, \tilde{\sigma}^2; 0.5, \beta_{\max})$$

where

$$\tilde{\sigma}^2 = \sigma^2 (\phi(j+i) + (1-\phi))$$

which will induce higher variance as we move from the top left corner of the P matrix, where the higher ranked block is, towards the bottom right corner, where the least preferred block is.

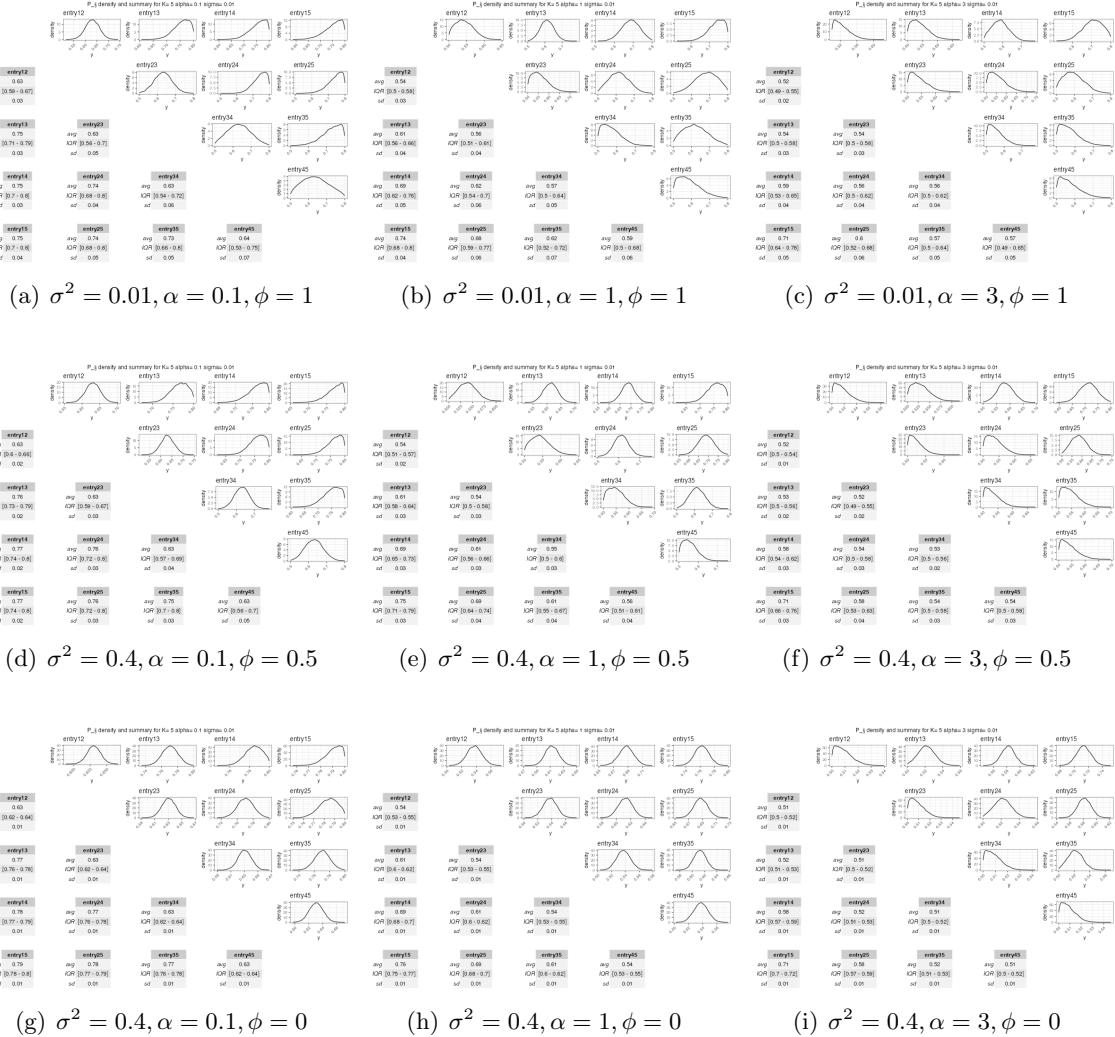


FIGURE 2. In the upper triangular side of each box above, we have the density of the p_{ij} , estimated over $n = 1000$ samples generated according to the combination of parameters σ^2, α, ϕ reported below each box. K is set to 5. In the lower triangular part of each box, we have tables with some summary statics, namely the mean, the interquartile range, and the standard deviation.

7. ESTIMATION

For the moment, we want to infer just $\theta = \{z, P, \alpha, \sigma^2\}$, meaning that we treat K as a known constant. We report below the posterior distribution that we want to estimate:

$$\begin{aligned}
 p(z, \alpha, \sigma^2, P \mid y) &= \frac{p(y \mid \alpha, \sigma^2, P) \cdot p(z \mid \gamma) \cdot p(\alpha) \cdot p(\sigma^2) \cdot p(P \mid \alpha, \sigma^2)}{\int p(y, \alpha, \sigma^2, P) dz d\alpha d\sigma^2 dP} \\
 &\propto p(y \mid \alpha, \sigma^2, P) \cdot p(z \mid \gamma) \cdot p(\alpha) \cdot p(\sigma^2) \cdot p(P \mid \alpha, \sigma^2) \\
 &= \prod_{i=2}^{N-1} \prod_{j=i}^N \text{Binomial}(y_{ij} \mid n_{ij}, p_{z_{ij}}) \cdot \text{DirichletMultinomial}(z \mid n, \gamma) \\
 &\quad \cdot \text{Unif}(\alpha \mid 0, 3) - \text{Unif}(\sigma^2 \mid 0, 1) \\
 &\quad \cdot \prod_{k=1}^{K-1} \text{TruncatedNormal}(L_{(k)} \mid \mu_k, \sigma^2, 0.5, \beta_{\max}) \\
 (18) \quad &= \prod_{i=2}^{N-1} \prod_{j=i}^N \binom{n_{ij}}{y_{ij}} (p_{z_{ij}})^{y_{ij}} (1 - p_{z_{ij}})^{n_{ij} - y_{ij}} \cdot \frac{\Gamma(\sum_{k=1}^K \gamma_k)}{\prod_{k=1}^K \Gamma(\gamma_k)} \frac{\prod_{k=1}^K \Gamma(n'_k + \gamma_k)}{\Gamma(\sum_{k=1}^K (n_k + \gamma_k))} \\
 &\quad \cdot \mathbb{I}_{0 \leq \alpha \leq 3} \cdot \frac{1}{3} \cdot \mathbb{I}_{0 \leq \sigma^2 \leq 1} \\
 &\quad \cdot \prod_{k=1}^{K-1} \prod_{\{i^*, j^*\}} \frac{1}{\sigma} \frac{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{p_{i^* j^*} - \mu_k}{\sigma} \right)^2}}{\int_{-\infty}^{\beta} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{t - \mu_k}{\sigma} \right)^2} dt - \int_{-\infty}^{0.5} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{t - \mu_k}{\sigma} \right)^2} dt} \mathbb{I}_{0.5 \leq p_{i^* j^*} \leq \beta_{\max}}
 \end{aligned}$$

where $\{i^*, j^*\}$ is the set of entries in P such that $j - i = k$, that is they belong to the diagonal k , and in turn the level set $L_{(k)}$. On the other hand, μ_k , the mean for each level set is equal to

$$\left[\frac{(\beta_{\max} - 0.5)}{2K^\alpha} \times (k^\alpha + (k+1)^\alpha) + \frac{1}{2} \right]$$

The estimation strategy for (??) is a Metropolis-within-Gibbs MCMC algorithm. MCMC techniques are designed precisely for scenarios like ours, where direct estimation of the posterior is infeasible or analytically intractable. By generating a sequence of samples through a carefully constructed Markov chain, MCMC allows us to navigate the complex parameter space in a data-driven manner.

Referencing Muller's (1991) work, we present here below an hybrid version of a classical Metropolis-within-Gibbs algorithm that has some self-tuning, or adaptive, features in the variance of the proposal distribution, the parameter τ_θ^2

We start by reviewing the MCMC step for z , and then we move on exploring the other continuous parameters in the remaining of this section.

Algorithm 1 Metropolis-within-Gibbs update for z

Given $\left(z^{(t), \alpha^{(t)}, \sigma^2(t), P^{(t)}}\right)$

1. Randomly sample an order

$$\pi \sim \text{Random permutation of } \{1, 2, \dots, n\}$$

for $i = 1$ to N do

1. Compute the difference $\{d_{k, z_{\pi(i)}^{(t)}}\}$ between each label $k = 1, \dots, K$ and the current state $z_{\pi(i)}^{(t)}$

$$d_{k, z_{\pi(i)}^{(t)}} = k - z_{\pi(i)}^{(t)}$$

2. Compute and normalize the probabilities

$$p_k = \frac{p\left(d_{k, z_{\pi(i)}^{(t)}}\right)}{\sum_{k=1}^K p\left(d_{k, z_{\pi(i)}^{(t)}}\right)} \quad \text{where} \quad d_{k, z_{\pi(i)}^{(t)}} \sim \text{Normal}\left(0, \tau_{z_{\pi(i)}}^2\right)$$

3. $z'_{\pi(i)} \leftarrow k'$ sampled from $\{k = 1, \dots, K\} \setminus \{z_{\pi(i)}^{(t)}\}$ with probability p_k

4. Take

$$z_{\pi(i)}^{(t+1)} = \begin{cases} z_{\pi(i)}^{(t)} & \text{with probability } 1 - r', \\ z'_{\pi(i)} & \text{with probability } r', \end{cases}$$

where

$$r' = \log(1) \wedge$$

$$(19) \quad \begin{aligned} & \log p\left(z'_{\pi(i)} \mid \alpha^{(t+1)}, \sigma^2(t+1), P^{(t)}, \left\{z_{\pi(j)}^{(t+1)} \mid j < i\right\}, \left\{z_{\pi(j)}^{(t)} \mid j > i\right\}\right) \\ & - \log p\left(z_{\pi(i)}^{(t)} \mid \left\{z_{\pi(j)}^{(t+1)} \mid j < i\right\}, \left\{z_{\pi(j)}^{(t)} \mid j > i\right\}, \alpha^{(t)}, \sigma^2(t), P^{(t)}, \right) \end{aligned}$$

end for

In algorithm (??) we have that (??) is equal to:

$$(20) \quad \sum_{i=2}^{N-1} \sum_{j=i}^N \text{LogBinomial}\left(y_{ij} \mid n_{ij}, p_{z'_i z'_j}^{(t)}\right) - \sum_{i=2}^{N-1} \sum_{j=i}^N \text{LogBinomial}\left(y_{ij} \mid n_{ij}, p_{z_i^{(t)} z_j^{(t)}}^{(t)}\right)$$

$$(21) \quad + \log \left(\frac{\Gamma(\sum_{k=1}^K \gamma_k)}{\prod_{k=1}^K \Gamma(\gamma_k)} \frac{\prod_{k=1}^K \Gamma(n'_k + \gamma_k)}{\Gamma(\sum_{k=1}^K (n'_k + \gamma_k))} \right) - \log \left(\frac{\Gamma(\sum_{k=1}^K \gamma_k)}{\prod_{k=1}^K \Gamma(\gamma_k)} \frac{\prod_{k=1}^K \Gamma(n_k^{(t)} + \gamma_k)}{\Gamma(\sum_{k=1}^K (n_k^{(t)} + \gamma_k))} \right)$$

since the joint prior on the level sets, that is

$$\sum_{k=1}^{K-1} \text{LogTruncatedNormal}\left(L_{(k)} \mid \left[\frac{(\beta_{\max} - 0.5)}{2K^{\alpha'}} \times (k^{\alpha'} + (k+1)^{\alpha'}) + \frac{1}{2}\right], \sigma^{2(t)}, 0.5, \beta_{\max}\right)$$

, the log prior probability on σ^2 , namely $\log p(\sigma)$, and the log prior probability on α are not affected by the new proposed value of z and therefore they subtract out.

This algorithm is meant to propose with higher probability those labels that are adjacent to the current one, and the Normal distribution centred on zero implies that the further the distance. Furthermore, the adaptive variance of the proposal $\tau_{z_{\pi(i)}}^2$ distribution controls the probability of labels that are more far apart from the current one.

Choosing a correct $\tau_{z_{\pi(i)}}^2$ value is not straightforward, and we choose to resort to an adaptive algorithm to elicitate a correct proposal variance. We proceed as in Roberts, Rosenthal 2012. For each of the i -th labels we create an associated variable ls_i giving the logarithm of the standard deviation to be used when proposing a normal increment to variable i . We begin with $ls_i = \log(0.04)$ for all i (corresponding to 0.2 proposal standard deviation). After the n -th "batch" of 50 iterations, we update each ls_i by adding or subtracting an adaption amount $\delta(n)$. The adapting attempts to make the acceptance rate of proposals for variable i as close as possible to 0.234, following the literature practice Chris Sherlock12009. Specifically, we increase ls_i by $\delta(n)$ if the fraction of acceptances of variable i was more than 0.234 on the n -th batch, or decrease ls_i by $\delta(n)$ if it was less.

Intuitively, if the acceptance rate for a particular label i is too low, we want the proposal to explore neighboring labels. Conversely, if the acceptance rate is too high, we aim to sample labels further away.

Now, let us investigate the MCMC move for α and the remaining parameters. The adaptive structure is maintained as for the parameter z , with the required adaptation to a continuous context. The fact that α , σ^2 and P are continuous parameters actually simplifies the notation and the expressions, since there is no need of the intermediate passage in which we compute the difference

In algorithm (??) we have (??) which is given by

$$\begin{aligned} & \sum_{k=1}^{K-1} \text{LogTruncatedNormal}\left(L_{(k)}^{(t)} \mid \left[\frac{(\beta_{\max} - 0.5)}{2K^{\alpha'}} \times (k^{\alpha'} + (k+1)^{\alpha'}) + \frac{1}{2}\right], (\sigma^2)^{(t)}, 0.5, \beta_{\max}\right) \\ & - \sum_{k=1}^{K-1} \text{LogTruncatedNormal}\left(L_{(k)}^{(t)} \mid \left[\frac{(\beta_{\max} - 0.5)}{2K^{\alpha^{(t)}}} \times (k^{\alpha^{(t)}} + (k+1)^{\alpha^{(t)}}) + \frac{1}{2}\right], (\sigma^2)^{(t)}, 0.5, \beta_{\max}\right) \\ (25) \quad & + \text{LogUnif}(\alpha' \mid 0, 3) - \text{LogUnif}((\sigma^2)^{(t)} \mid 0, 1) \end{aligned}$$

since the log-likelihood $\log(p(y \mid z, P))$, the log prior probability on z , namely $\log p(z \mid \gamma)$, and the log prior prior probability on σ^2 , namely $\log p(\sigma)$, are not affected by the new proposed value of α and therefore they subtract out.

Algorithm 2 Metropolis-within-Gibbs update for α

Given $(z^{(t+1)}, P^{(t)}, \sigma^{2(t)})$

1. Sample

$$(22) \quad \alpha' \text{from Normal} \left(\alpha^{(t-1)}, (\tau_\alpha^2)^{(t-1)} \right)$$

2. Take

$$(23) \quad \alpha^{(t+1)} = \begin{cases} \alpha^{(t)} & \text{with probability } 1 - r', \\ \alpha' & \text{with probability } r', \end{cases}$$

where

$$(24) \quad r' = \log(1) \wedge \log p \left(\alpha' | z^{(t+1)}, \alpha^{(t)}, P^{(t)}, \sigma^{2(t)} \right) - \log p \left(\alpha^{(t)} | z^{(t+1)}, \alpha^{(t)}, P^{(t)}, \sigma^{2(t)} \right)$$

Algorithm 3 Metropolis-within-Gibbs update for σ^2

Given $(z^{(t+1)}, \alpha^{(t+1)}, P^{(t)})$

1. Sample

$$(26) \quad (\sigma^2)' \text{from Normal} \left((\sigma^2)^{(t-1)}, (\tau_{\sigma^2}^2)^{(t-1)} \right)$$

2. Take

$$(27) \quad (\sigma^2)^{(t+1)} = \begin{cases} (\sigma^2)^{(t)} & \text{with probability } 1 - r', \\ (\sigma^2)' & \text{with probability } r', \end{cases}$$

where

$$(28) \quad \log p \left((\sigma^2)' | z^{(t+1)}, \alpha^{(t+1)}, (\sigma^2)^{(t)}, P^{(t)} \right) - \log p \left((\sigma^2)^{(t)} | z^{(t+1)}, \alpha^{(t+1)}, (\sigma^2)^{(t)}, P^{(t)}, \right)$$

In algorithm (??) we have that (??) is given by

$$(29) \quad \begin{aligned} & \sum_{k=1}^{K-1} \text{LogTruncatedNormal} \left(L_{(k)}^{(t)} | \left[\frac{(\beta_{\max} - 0.5)}{2K^{\alpha^{(t)}}} \times \left(k^{\alpha^{(t)}} + (k+1)^{\alpha^{(t)}} \right) + \frac{1}{2} \right], (\sigma^2)', 0.5, \beta_{\max} \right) \\ & - \sum_{k=1}^{K-1} \text{LogTruncatedNormal} \left(L_{(k)}^{(t)} | \left[\frac{(\beta_{\max} - 0.5)}{2K^{\alpha^{(t)}}} \times \left(k^{\alpha^{(t)}} + (k+1)^{\alpha^{(t)}} \right) + \frac{1}{2} \right], (\sigma^2)^{(t)}, 0.5, \beta_{\max} \right) \\ & + \text{LogUnif} \left(\alpha^{(t)} | 0, 3 \right) - \text{LogUnif} \left((\sigma^2)^{(t)} | 0, 1 \right) \end{aligned}$$

since the log-likelihood $\log(p(y | z, P))$, the log prior probability on z , namely $\log p(z | \gamma)$, and the log prior prior probability on α^2 , namely $\log p(\alpha)$, are not affected by the new proposed value of α and therefore they subtract out.

Algorithm 4 Metropolis-within-Gibbs update for P

Given $(z^{(t+1)}, \alpha^{(t+1)}, \sigma^{2(t+1)}, P^{(t)})$

for $i = 1$ to $K - 1$ **do**

for $j = i + 1$ to K **do**

1. Sample

$$(30) \quad p'_{ij} \text{from Normal} \left(p_{ij}^{(t)}, (\tau_{p_{ij}}^2)^{(t-1)} \right)$$

2. Take

$$(31) \quad p_{ij}^{(t+1)} = \begin{cases} p_{ij}^{(t)} & \text{with probability } 1 - r', \\ p'_{ij} & \text{with probability } r', \end{cases}$$

where

$$r' = \log(1) \wedge$$

$$(32) \quad \begin{aligned} & \log p \left(p'_{i,j} \mid z^{(t+1)}, \alpha^{(t+1)}, \sigma^{2(t+1)}, \left\{ p_{i^*,j^*}^{(t+1)} \mid i^* < i, j^* < j \right\}, \left\{ p_{i^*,j^*}^{(t)} \mid i^* > i, j^* > j \right\} \right) \\ & - \log p \left(p_{i,j}^{(t)} \mid z^{(t+1)}, \alpha^{(t+1)}, \sigma^{2(t+1)}, \left\{ p_{i^*,j^*}^{(t+1)} \mid i^* < i, j^* < j \right\}, \left\{ p_{i^*,j^*}^{(t)} \mid i^* > i, j^* > j \right\} \right) \end{aligned}$$

end for

end for

In algorithm (??), (??) is given by:

$$(33) \quad \begin{aligned} & \sum_{i=2}^{N-1} \sum_{j=i}^N \text{LogBinomial} \left(y_{ij} \mid n_{ij}, p'_{z_i^{(t)} z_j^{(t)}} \right) - \sum_{i=2}^{N-1} \sum_{j=i}^N \text{LogBinomial} \left(y_{ij} \mid n_{ij}, p_{z_i^{(t)} z_j^{(t)}}^{(t)} \right) \\ & \sum_{k=1}^{K-1} \text{LogTruncatedNormal} \left(L'_{(k)} \mid \left[\frac{(\beta_{\max} - 0.5)}{2K^{\alpha^{(t)}}} \times \left(k^{\alpha^{(t)}} + (k+1)^{\alpha^{(t)}} \right) + \frac{1}{2} \right], (\sigma^2)^{(t)}, 0.5, \beta_{\max} \right) \\ & - \sum_{k=1}^{K-1} \text{LogTruncatedNormal} \left(L_{(k)}^{(t)} \mid \left[\frac{(\beta_{\max} - 0.5)}{2K^{\alpha^{(t)}}} \times \left(k^{\alpha^{(t)}} + (k+1)^{\alpha^{(t)}} \right) + \frac{1}{2} \right], (\sigma^2)^{(t)}, 0.5, \beta_{\max} \right) \end{aligned}$$

As before, the adaptive proposal variance is specific to each upper-triangular entry of P . Therefore, each entry's acceptance rate is monitored and the variance is adjusted accordingly.m

→ Insert here plots of convergence to the acceptance ratio

8. POINT ESTIMATE, MODEL SELECTION, AND INFERENCE

While algorithmic methods produce a single estimated partition, our model offers the entire posterior distribution across different node partitions. We are comparing the results from the simulation study via the following measures

We obtain the point estimate \hat{z} from the MCMC samples in two ways. The first one is via the MAP estimate, meaning that partition which maximises the a posteriori distribution.

$$\begin{aligned}
 \hat{z}_{\text{MAP}}(y) &= \arg \max_z f(z | y) \\
 &= \arg \max_z \frac{f(y | z), g(z)}{\int_Z f(x | z), g(z) dz} \\
 (35) \quad &= \arg \max_z f(x | z) g(z)
 \end{aligned}$$

The second one is the partition that minimises the lowest averaged variation of information (VI distance) from the other clusterings, denoted in the following as \hat{z}_{lbVI} .

The VI fully utilise this posterior and it founded on a decision-theoretic approach introduced by Wade and Ghahramani (2018) for block modelling. This involves summarizing posterior distributions using the variation of information (VI) metric, developed by Meilă (2007), which measures the distance between two clusterings by comparing their individual and joint entropies. The VI metric ranges from 0 to $\log 2N$, where N represents the number of nodes. Intuitively, the VI metric quantifies the amount of information contained in two clusterings relative to the shared information between them. As a result, it decreases towards 0 as the overlap between two partitions increases. The variation of information is a true distance since it obeys the triangle inequality.

Suppose we have two partitions of a set X and Y of a set A into disjoint subsets, namely $X = \{X_1, X_2, \dots, X_k\}$ and $Y = \{Y_1, Y_2, \dots, Y_l\}$.

Let: $n = \sum_i |X_i| = \sum_j |Y_j| = |A|$ $p_i = |X_i|/n$ and $q_j = |Y_j|/n$ $r_{ij} = |X_i \cap Y_j|/n$

Then the variation of information between the two partitions is:

$$\text{VI}(X; Y) = - \sum_{i,j} r_{ij} [\log(r_{ij}/p_i) + \log(r_{ij}/q_j)]$$

To compare different models we use the WAIC loss, which yields practical and theoretical advantages with respect to other losses and has direct connections with Bayesian leave-one-out cross-validation, thus providing a measure of edge predictive accuracy.

Moreover, the calculation of the WAIC only requires posterior samples of the log-likelihoods for the edges: $\log p(y_{ij}|z, P, \alpha) = y_{ij} \log p_{z_i, z_j} + (n_{ij} - y_{ij}) \log(1 - p_{z_i, z_j})$, $i = 2, \dots, N, j = 1, \dots, i - 1$. These quantities are already available to the user, since the MCMC chain provides sample both of z and P . The use of the WAIC in a context similar to the present one is documented in Durante and Legramanti (2020).

9. SIMULATION STUDY FROM THE UNORDERED MODEL N=100

In order to evaluate how well our model performs in a situation similar to our intended use, and measure its advantages compared to the best existing alternatives, we generated three simulated tournaments with 100 players from the Unordered Model.

We want to compare the quality of the Unordered model in recovering the ground truth from some data generated from the Unordered model itself. We compare its recovery performance with the one of the POMM model, which in the present context qualifies as a sort of 'mis-specified model', since the blocks do not present an inherent ordering.

Below, in figure (??), you may see the adjacency matrix of the simulated data. The darker is the pixel (i, j) , the more are the victories of player i vs player j .

On the side, the different colours testify the different block membership of each single player.

TABLE 1. Time performance

Fitted Model	Seconds per iteration			Expected time for 30000 iterations		
	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0.013 sec	0.015 sec	0.018 sec	6.5 min	7.5 min	9 min
Unordered model	0.013 sec	0.015 sec	0.018 sec	6.5 min	7.5 min	9 min

For the estimation purposes, we run 4 different chains of 30000 iterations each. We report in table (??) the execution time for the MCMC. Within that table, and in all other tables that will follow in the simulation study column (a) refers to the case $K = 3$, column (b) refers to the case $K = 5$, and column (c) to the case $K = 9$.

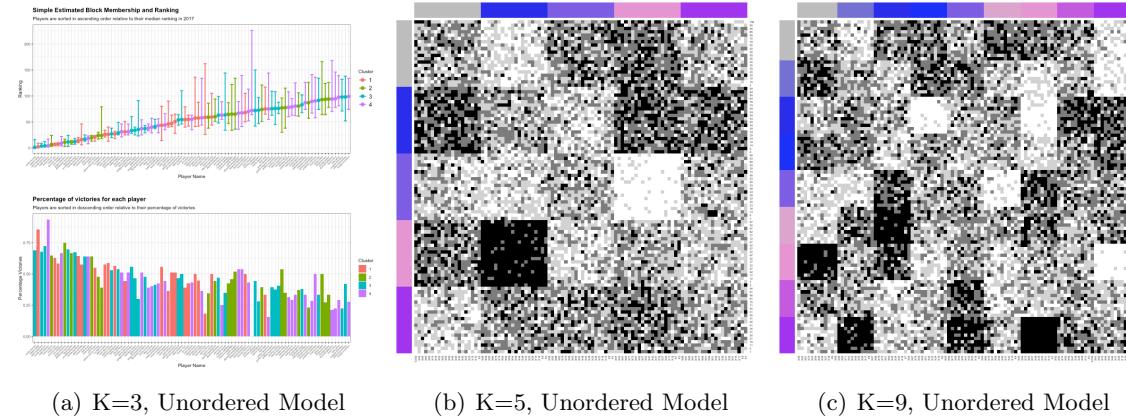


FIGURE 3. Adjacency Matrices simulated via the Unordered Model

Each chain is initiated with different starting values and different seeds. The initiation values are saved in order to guarantee the reproducibility of the results.

For the POMM model, we need to choose an appropriate value for β_{\max} , which controls the maximum attainable value within the matrix P . Here we fix it arbitrarily at 0.85.

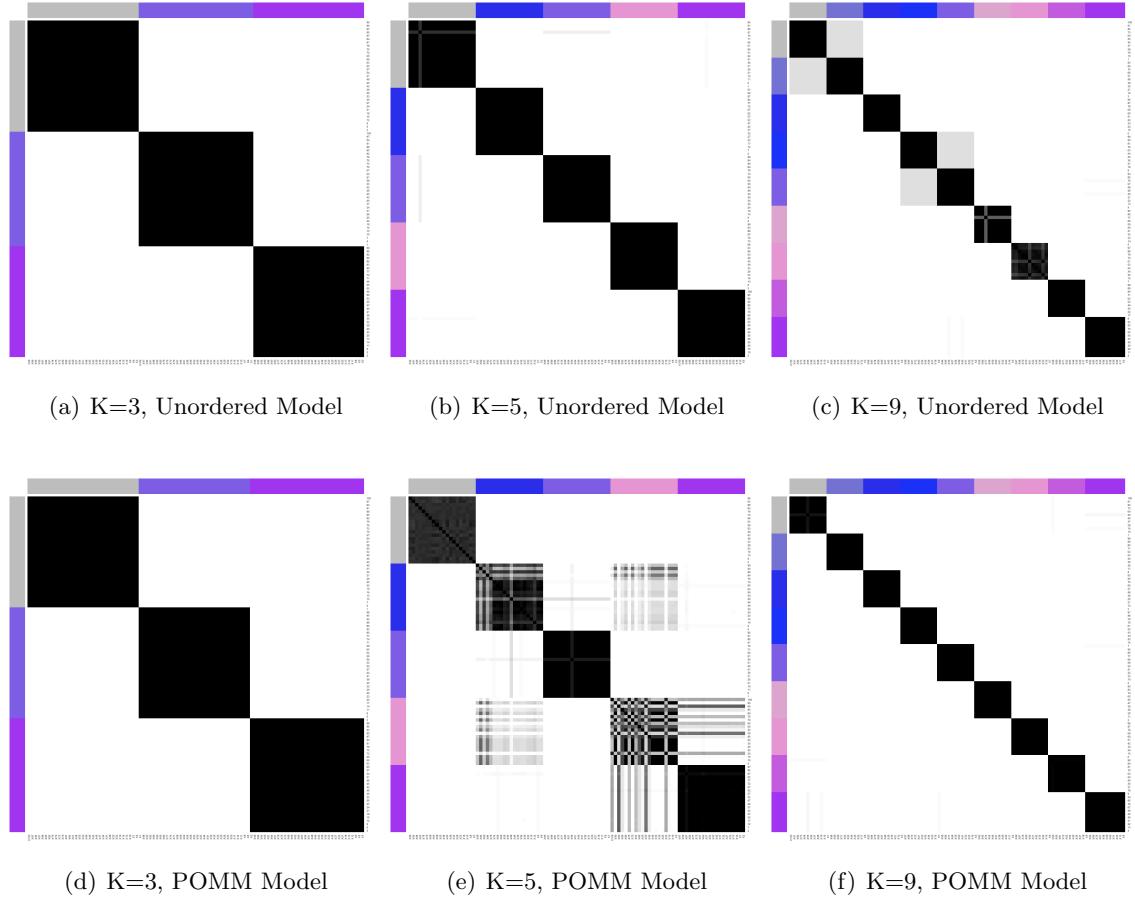


FIGURE 4. Co-Clustering Matrices obtained via the Unordered Model (above) and the POMM model (below).

The first results' plot we report is the one in figure (??). Each box contains a co-clustering matrix. Each pixel (i, j) represents the probability that two individuals are placed within the same cluster. The darker the pixel, the higher the probability. Colours on the side signal the true membership of each player.

In the first row, containing figure ??, figure ??, and figure ??, we have the co-clustering matrix for the Simple model estimated on the data generated according to the Simple model itself. This means that the blocks have no inherent ordering, and the model here is correctly specified. We can notice a very good recovery of the true membership.

In the second row instead, the one which contains figure ??, figure ??, and ??, we may observe the co-clustering matrix for the POMM model estimated on the data generated via

the Simple one. Therefore, the model is misspecified, but we may notice that the recovery performance is quite competitive with the Simple one.

TABLE 2. P summary table
True Model Unordered, $N = 100$

Fitted Model	MAE			% within-95% CI interval			CI interval length		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0.07	0.15	0.15	66.67	40	25.00	0.07	0.12	0.09
Unordered model	0.01	0.13	0.13	100.00	40	41.67	0.02	0.05	0.19

In table (??), we report the summary table of the estimates for the P matrix, both obtained via the Unordered model and the POMM model. For each value of K (again, reported in the sub-columns (a), (b) and (c), respectively), we report the mean absolute error, meaning the mean absolute error $MAE = \frac{1}{(K \cdot (K-1))/2} \sum_{i=1}^{K-1} \sum_{j=i+1}^K (\hat{p}_{ij} - p_{ij}^*)$ where p_{ij}^* is the true value of that particular entry. Then we also report the percentage of the upper triangular entries of P that are contained by the estimated 95% credible intervals, obtained by computing the 95% higher posterior density region, and the average 95% credible interval length.

Here below, we report the ground truth of the P matrix for each case $K = 3, 5, 9$. Then, next and below we report the estimates, both for the POMM and the Unordered model, within the $\hat{P}_{\text{POMM}}^{K=k}$, $\hat{P}_{\text{Unordered}}^{K=k}$, respectively.

$$\begin{aligned}
P_{true}^{K=3} &= \begin{bmatrix} 0.500 & 0.230 & 0.631 \\ 0.770 & 0.500 & 0.327 \\ 0.369 & 0.673 & 0.500 \end{bmatrix} & \hat{P}_{\text{POMM}}^{K=3} &= \begin{bmatrix} 0.500 & 0.222 & 0.50 \\ 0.778 & 0.500 & 0.41 \\ 0.500 & 0.590 & 0.50 \end{bmatrix} \\
\hat{P}_{\text{Unordered}}^{K=3} &= \begin{bmatrix} 0.500 & 0.222 & 0.639 \\ 0.778 & 0.500 & 0.321 \\ 0.361 & 0.679 & 0.500 \end{bmatrix} \\
P_{true}^{K=5} &= \begin{bmatrix} 0.500 & 0.230 & 0.631 & 0.706 & 0.422 \\ 0.770 & 0.500 & 0.327 & 0.752 & 0.714 \\ 0.369 & 0.248 & 0.500 & 0.036 & 0.441 \\ 0.673 & 0.964 & 0.286 & 0.500 & 0.365 \\ 0.294 & 0.578 & 0.559 & 0.635 & 0.500 \end{bmatrix} & \hat{P}_{\text{POMM}}^{K=5} &= \begin{bmatrix} 0.500 & 0.228 & 0.614 & 0.438 & 0.532 \\ 0.772 & 0.500 & 0.515 & 0.473 & 0.562 \\ 0.386 & 0.485 & 0.500 & 0.403 & 0.421 \\ 0.562 & 0.527 & 0.597 & 0.500 & 0.467 \\ 0.468 & 0.438 & 0.579 & 0.533 & 0.500 \end{bmatrix} \\
\hat{P}_{\text{Unordered}}^{K=5} &= \begin{bmatrix} 0.500 & 0.226 & 0.622 & 0.522 & 0.560 \\ 0.774 & 0.500 & 0.511 & 0.427 & 0.571 \\ 0.378 & 0.439 & 0.500 & 0.372 & 0.435 \\ 0.478 & 0.591 & 0.628 & 0.500 & 0.350 \\ 0.440 & 0.429 & 0.565 & 0.650 & 0.500 \end{bmatrix} \\
P_{true}^{K=9} &= \begin{bmatrix} 0.500 & 0.230 & 0.631 & 0.706 & 0.422 & 0.765 & 0.720 & 0.554 & 0.231 \\ 0.770 & 0.500 & 0.327 & 0.752 & 0.714 & 0.363 & 0.197 & 0.512 & 0.118 \\ 0.369 & 0.559 & 0.500 & 0.036 & 0.441 & 0.542 & 0.034 & 0.795 & 0.770 \\ 0.673 & 0.635 & 0.280 & 0.500 & 0.365 & 0.458 & 0.262 & 0.525 & 0.722 \\ 0.294 & 0.235 & 0.803 & 0.446 & 0.500 & 0.082 & 0.764 & 0.567 & 0.553 \\ 0.248 & 0.637 & 0.966 & 0.488 & 0.565 & 0.500 & 0.712 & 0.435 & 0.636 \\ 0.964 & 0.458 & 0.738 & 0.205 & 0.525 & 0.230 & 0.500 & 0.475 & 0.020 \\ 0.578 & 0.542 & 0.236 & 0.475 & 0.769 & 0.278 & 0.364 & 0.500 & 0.382 \\ 0.286 & 0.918 & 0.288 & 0.433 & 0.882 & 0.447 & 0.980 & 0.618 & 0.500 \end{bmatrix} \\
\hat{P}_{\text{POMM}}^{K=9} &= \begin{bmatrix} 0.500 & 0.293 & 0.559 & 0.488 & 0.518 & 0.499 & 0.404 & 0.485 & 0.453 \\ 0.707 & 0.500 & 0.503 & 0.496 & 0.615 & 0.340 & 0.466 & 0.504 & 0.202 \\ 0.441 & 0.497 & 0.500 & 0.385 & 0.316 & 0.302 & 0.203 & 0.638 & 0.498 \\ 0.512 & 0.504 & 0.615 & 0.500 & 0.509 & 0.501 & 0.523 & 0.549 & 0.632 \\ 0.482 & 0.385 & 0.684 & 0.491 & 0.500 & 0.264 & 0.631 & 0.503 & 0.332 \\ 0.501 & 0.660 & 0.698 & 0.499 & 0.736 & 0.500 & 0.714 & 0.586 & 0.532 \\ 0.596 & 0.534 & 0.797 & 0.477 & 0.369 & 0.286 & 0.500 & 0.535 & 0.201 \\ 0.515 & 0.496 & 0.362 & 0.451 & 0.497 & 0.414 & 0.465 & 0.500 & 0.358 \\ 0.547 & 0.798 & 0.502 & 0.368 & 0.668 & 0.468 & 0.799 & 0.642 & 0.500 \end{bmatrix}
\end{aligned}$$

$$\hat{P}_{\text{Unordered}}^{K=9} = \begin{bmatrix} 0.500 & 0.351 & 0.620 & 0.484 & 0.576 & 0.756 & 0.371 & 0.449 & 0.456 \\ 0.649 & 0.500 & 0.445 & 0.466 & 0.713 & 0.371 & 0.432 & 0.519 & 0.210 \\ 0.380 & 0.555 & 0.500 & 0.458 & 0.325 & 0.302 & 0.149 & 0.761 & 0.739 \\ 0.516 & 0.534 & 0.542 & 0.500 & 0.481 & 0.474 & 0.582 & 0.536 & 0.652 \\ 0.424 & 0.287 & 0.675 & 0.519 & 0.500 & 0.290 & 0.616 & 0.401 & 0.375 \\ 0.244 & 0.629 & 0.698 & 0.526 & 0.710 & 0.500 & 0.717 & 0.595 & 0.568 \\ 0.629 & 0.568 & 0.851 & 0.418 & 0.384 & 0.283 & 0.500 & 0.559 & 0.058 \\ 0.551 & 0.481 & 0.239 & 0.464 & 0.599 & 0.405 & 0.441 & 0.500 & 0.323 \\ 0.544 & 0.790 & 0.261 & 0.348 & 0.625 & 0.432 & 0.942 & 0.677 & 0.500 \end{bmatrix}$$

TABLE 3. z summary table
True Model Unordered, $N = 100$

Method	VI distance _{MAP}			VI distance _{VI lb}			WAIC		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0	0	0	0	0.56	0	48883.71 145.40	53337.62 152.84	51244.98 153.63
Unordered model	0	0	0	0	0.00	0	46521.89 144.55	52382.24 153.98	49144.47 164.65

In table (??), we report some summary statistics for the parameter z . As above, we have columns (a),(b) and (c) representing the case $K = 3, 5, 9$ respectively.

The first indicator is the $VI\text{distance}_{\text{MAP}}$ computed between the true partition z^* and the point estimate \hat{z}^{MAP} obtained with the maximum a posteriori estimate (MAP).

The second one is $VI\text{distance}_{\text{VI lb}}$ computed between the true partition z^* and the point estimate $\hat{z}^{\text{VI lb}}$ obtained with the partition attaining the VI lower bound.

The third one is the WAIC estimate, along with its standard error below. We see that the Unordered Model is the one preferred in this case.

TABLE 4. POMM hyperparameters summary table
True Model Unordered, $N = 100$

Fitted Model	$\hat{\theta}$			95% CI		
	(a)	(b)	(c)	(a)	(b)	(c)
σ	0.51	0.57	0.64	[0.13 0.89]	[0.24 0.9]	[0.34 0.9]
α	0.48	0.52	0.59	[0.12 0.87]	[0.15 0.88]	[0.19 0.9]

In table (??), we report the results for the hyperparameters of the *POMM* model. In the $\hat{\theta}$ column we report the estimates both for α and σ , while on the right we have their 95% Credible Interval. Given that the data were generated according the Unordered model, we do not have a ground truth to which these results should be compared. However, we can

still try to make sense of these values by inspecting the properties of the induced P^{POMM} matrix resulting from the estimates.

The Unordered model has a prior over $P^{\text{Unordered}} \sim \text{Beta}(1, 1)$, then we may expect P^{POMM} itself to get as closer as needed to a uniform distribution by selecting the appropriate combination of (α, σ) .

9.0.1. *Unordered Model check.* In this subsection, we report some diagnostic checks for the algorithm of the Markov Chains, to assess convergence, quality of mixing, and the overall behaviour of the Metropolis-within-Gibbs algorithm.

TABLE 5. z diagnostic table
True Model Unordered, $N = 100$

Fitted Model	\overline{ESS}			\overline{ACF}_{30}			$\overline{\%accepted}$			$\overline{Gelman - Rubin}$		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	7558.60	301.04	2880.01	0	0.01	0.08	0.02	0.04	0.15	1.61	1.29	1.31
Unordered	6544.85	253.85	140.78	0	0.03	0.11	0.01	0.56	0.15	1.04	1.23	1.26

In table (??) we report the diagnostics for the z parameter. Since the parameter is a label vector with 100 entries in this case, we compute the relevant statistics for each single entries an then we report the average.

- The first statistics is the Effective Sample Size (ESS) averaged over individuals $i = 1, \dots, N$, which denotes a fairly good sample size. The average is taken as follows:

$$\overline{ESS} = \frac{1}{n} \sum_{i=1}^N ESS_i$$

The same is applied also to the other diagnostic metrics.

- Then, we report the average autocorrelation, \overline{ACF}_{30} , computed with a lag of 30 iterations. This values are close to zero, meaning that there is very little correlation within the chain.
- In the third column, we report the average acceptance rate $\overline{\%accepted}$. These are significantly lower than the target acceptance rate that should be hit by the adaptive MCMC, which is 22%. However, the estimates are capable of correctly recovering the true partition. Combining the two facts, we may hypothesise that the simulation study has too "many" data, and proposing for a given individual i , who has already been assigned to the true block, a block different from the true one, leads to a drop in the likelihood which is too large, and as a consequence, we always reject the other labels.
- Finally, we compute the median Gelman-Rubin statistics for each entry, $\overline{Gelman - Rubin}$. Gelman and Rubin (1992) propose a general approach to monitoring convergence of MCMC output in which $m > 1$ parallel chains are run with starting values that are overdispersed relative to the posterior distribution. Convergence is diagnosed when the chains have 'forgotten' their initial values, and the output from all chains is indistinguishable. The Gelman-Rubin diagnostic is applied to a single variable

from the chain. It is based a comparison of within-chain and between-chain variances, and is similar to a classical analysis of variance. Values substantially above 1 indicate lack of convergence. If the chains have not converged,

TABLE 6. P diagnostic table
True Model Unordered, $N = 100$

Fitted Model	\overline{ESS}			$\overline{ACF_{30}}$			$\overline{\% \text{accepted}}$			$Gelman - Rubin$		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	345.67	639.8	704.53	0.35	0.05	0.06	7.53	19.06	21.60	1.03	4.23	1.01
Unordered	570.00	828.1	937.67	0.00	0.04	0.03	11.66	18.40	29.06	1.00	1.01	2.23

In table (??) we report the same diagnostics checks for the z parameter. The only difference is that here we do not average the diagnostics indicators over the individuals $i = 1, \dots, N$, but instead over the upper-triangular P indices: $\{i = 1, \dots, K - 1, j = i + 1, \dots, K\}$.

Just as an example, the average ESS in this case is obtained as

$$\overline{ESS} = \frac{1}{(K \cdot (K - 1))/2} \sum_{i=1}^{N-1} \sum_{j=i+1}^N ESS_{i,j}$$

TABLE 7. POMM hyperparameters diagnostic table
True Model Unordered, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
σ	981	1041	842	0.01	0.01	0.01	39.03	36.59	32	1.08	1	1.02
α	26	17	19	0.82	0.82	0.82	25.36	25.18	24.74	1.33	1.21	1.19

Finally, in table (??) we report again the same diagnostics, but since both α and σ are one-dimensional, we are presenting the diagnostics itself, without any average.

10. SIMULATION STUDY FROM THE WST MODEL N=100

Differing from the Unordered model, the WST model has the constraint that the upper triangular entries of the matrix P are all greater than 0.5.

We want to compare the quality of the POMM model in recovering the ground truth from some data generated from the WST model itself. We compare its recovery performance with the one of the POMM model, which in the present context qualifies as a sort of 'mis-specified model', since the blocks do not present an inherent ordering.

Below, in figure (??), you may see the adjacency matrix of the simulated data. The darker is the pixel (i, j) , the more are the victories of player i vs player j .

On the side, the different colours testify the different block membership of each single player.

TABLE 8. Time performance

Fitted Model	Seconds per iteration			Expected time for 30000 iterations		
	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0.013 sec	0.015 sec	0.018 sec	6.5 min	7.5 min	9 min
WST model	0.013 sec	0.015 sec	0.018 sec	6.5 min	7.5 min	9 min

For the estimation purposes, we run 4 different chains of 30000 iterations each. We report in table (??) the execution time for the MCMC. Within that table, and in all other tables that will follow in the simulation study column (a) refers to the case $K = 3$, column (b) refers to the case $K = 5$, and column (c) to the case $K = 9$.

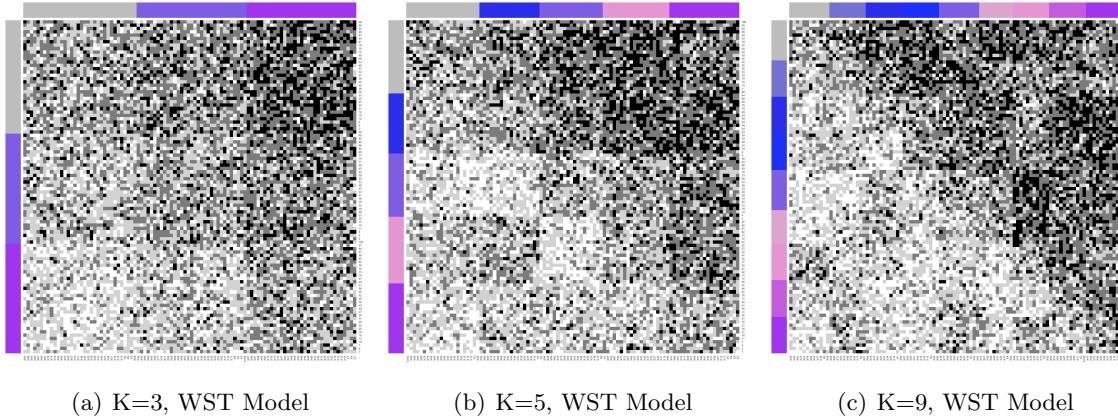


FIGURE 5. Adjacency Matrices simulated via the WST Model

Each chain is initiated with different starting values and different seeds. The initiation values are saved in order to guarantee the reproducibility of the results.

For the POMM model, we need to choose an appropriate value for β_{\max} , which controls the maximum attainable value within the matrix P . Here we fix it arbitrarily at 0.85.

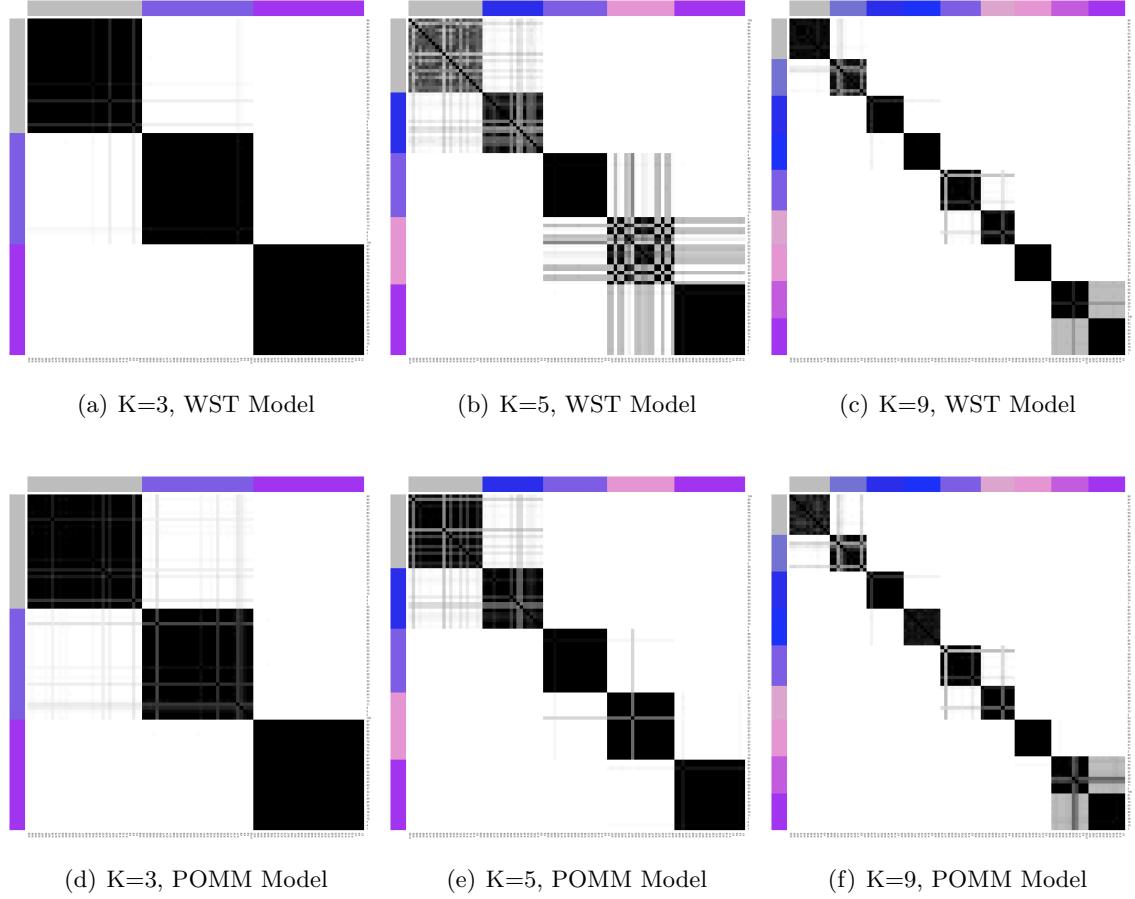


FIGURE 6. Co-Clustering Matrices obtained via the WST Model (above) and the POMM model (below).

The first results' plot we report is the one in figure (??). Each box contains a co-clustering matrix. Each pixel (i, j) represents the probability that two individuals are placed within the same cluster. The darker the pixel, the higher the probability. Colours on the side signal the true membership of each player.

In the first row, containing figure ??, figure ??, and figure ??, we have the co-clustering matrix for the WST model estimated on the data generated according to the WST model itself. This means that the blocks have no inherent ordering, and the model here is correctly specified. We can notice a very good recovery of the true membership.

In the second row instead, the one which contains figure ??, figure ??, and ??, we may observe the co-clustering matrix for the POMM model estimated on the data generated via

the WST one. Therefore, the model is misspecified, but we may notice that the recovery performance is quite competitive with the WST one.

TABLE 9. P summary table
True Model WST, $N = 100$

Fitted Model	<i>MAE</i>			% within-95% CI interval			CI interval length		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0.02	0.06	0.06	100	40	72.22	0.12	0.04	0.16
WST model	0.00	0.06	0.06	100	50	69.44	0.02	0.10	0.15

In table (??), we report the summary table of the estimates for the P matrix, both obtained via the WST model and the POMM model. For each value of K (again, reported in the sub-columns (a), (b) and (c), respectively), we report the mean absolute error, meaning the mean absolute error $MAE = \frac{1}{(K \cdot (K-1))/2} \sum_{i=1}^{K-1} \sum_{j=i+1}^K (\hat{p}_{ij} - p_{ij}^*)$ where p_{ij}^* is the true value of that particular entry. Then we also report the percentage of the upper triangular entries of P that are contained by the estimated 95% credible intervals, obtained by computing the 95% higher posterior density region, and the average 95% credible interval length.

Here below, we report the ground truth of the P matrix for each case $K = 3, 5, 9$. Then, next and below we report the estimates, both for the POMM and the WST model, within the $\hat{P}_{\text{POMM}}^{K=k}$, $\hat{P}_{\text{WST}}^{K=k}$, respectively.

$$P_{true}^{K=3} = \begin{bmatrix} 0.500 & 0.586 & 0.736 \\ 0.414 & 0.500 & 0.623 \\ 0.264 & 0.377 & 0.500 \end{bmatrix} \quad \hat{P}_{\text{POMM}}^{K=3} = \begin{bmatrix} 0.500 & 0.600 & 0.754 \\ 0.400 & 0.500 & 0.622 \\ 0.246 & 0.378 & 0.500 \end{bmatrix}$$

$$\hat{P}_{\text{WST}}^{K=3} = \begin{bmatrix} 0.500 & 0.592 & 0.739 \\ 0.408 & 0.500 & 0.627 \\ 0.261 & 0.373 & 0.500 \end{bmatrix}$$

$$P_{true}^{K=5} = \begin{bmatrix} 0.500 & 0.586 & 0.736 & 0.765 & 0.658 \\ 0.414 & 0.500 & 0.623 & 0.782 & 0.768 \\ 0.264 & 0.218 & 0.500 & 0.514 & 0.665 \\ 0.377 & 0.486 & 0.232 & 0.500 & 0.637 \\ 0.235 & 0.342 & 0.335 & 0.363 & 0.500 \end{bmatrix} \quad \hat{P}_{\text{POMM}}^{K=5} = \begin{bmatrix} 0.500 & 0.589 & 0.731 & 0.687 & 0.718 \\ 0.411 & 0.500 & 0.703 & 0.640 & 0.700 \\ 0.269 & 0.297 & 0.500 & 0.643 & 0.658 \\ 0.313 & 0.360 & 0.357 & 0.500 & 0.653 \\ 0.282 & 0.300 & 0.342 & 0.347 & 0.500 \end{bmatrix}$$

$$\hat{P}_{\text{WST}}^{K=5} = \begin{bmatrix} 0.500 & 0.566 & 0.660 & 0.700 & 0.710 \\ 0.434 & 0.500 & 0.646 & 0.675 & 0.703 \\ 0.340 & 0.354 & 0.500 & 0.666 & 0.665 \\ 0.300 & 0.325 & 0.334 & 0.500 & 0.652 \\ 0.290 & 0.297 & 0.335 & 0.348 & 0.500 \end{bmatrix}$$

$$P_{true}^{K=9} = \begin{bmatrix} 0.500 & 0.586 & 0.736 & 0.765 & 0.658 & 0.787 & 0.770 & 0.708 & 0.587 \\ 0.414 & 0.500 & 0.623 & 0.782 & 0.768 & 0.636 & 0.574 & 0.692 & 0.544 \\ 0.264 & 0.335 & 0.500 & 0.514 & 0.665 & 0.703 & 0.513 & 0.798 & 0.789 \\ 0.377 & 0.363 & 0.230 & 0.500 & 0.637 & 0.672 & 0.598 & 0.697 & 0.771 \\ 0.235 & 0.213 & 0.426 & 0.292 & 0.500 & 0.531 & 0.786 & 0.713 & 0.707 \\ 0.218 & 0.364 & 0.487 & 0.308 & 0.337 & 0.500 & 0.767 & 0.663 & 0.739 \\ 0.486 & 0.297 & 0.402 & 0.202 & 0.322 & 0.211 & 0.500 & 0.678 & 0.507 \\ 0.342 & 0.328 & 0.214 & 0.303 & 0.413 & 0.229 & 0.261 & 0.500 & 0.643 \\ 0.232 & 0.469 & 0.233 & 0.287 & 0.456 & 0.293 & 0.493 & 0.357 & 0.500 \end{bmatrix}$$

$$\hat{P}_{\text{POMM}}^{K=9} = \begin{bmatrix} 0.500 & 0.572 & 0.666 & 0.707 & 0.723 & 0.712 & 0.705 & 0.626 & 0.695 \\ 0.428 & 0.500 & 0.646 & 0.689 & 0.716 & 0.725 & 0.667 & 0.656 & 0.635 \\ 0.334 & 0.354 & 0.500 & 0.645 & 0.649 & 0.656 & 0.614 & 0.701 & 0.729 \\ 0.293 & 0.311 & 0.355 & 0.500 & 0.651 & 0.666 & 0.666 & 0.687 & 0.751 \\ 0.277 & 0.284 & 0.351 & 0.349 & 0.500 & 0.631 & 0.714 & 0.669 & 0.692 \\ 0.288 & 0.275 & 0.344 & 0.334 & 0.369 & 0.500 & 0.667 & 0.732 & 0.674 \\ 0.295 & 0.333 & 0.386 & 0.334 & 0.286 & 0.333 & 0.500 & 0.687 & 0.668 \\ 0.374 & 0.344 & 0.299 & 0.313 & 0.331 & 0.268 & 0.313 & 0.500 & 0.603 \\ 0.305 & 0.365 & 0.271 & 0.249 & 0.308 & 0.326 & 0.332 & 0.397 & 0.500 \end{bmatrix}$$

$$\hat{P}_{\text{WST}}^{K=9} = \begin{bmatrix} 0.500 & 0.601 & 0.675 & 0.711 & 0.702 & 0.744 & 0.702 & 0.659 & 0.709 \\ 0.399 & 0.500 & 0.623 & 0.671 & 0.734 & 0.702 & 0.653 & 0.647 & 0.581 \\ 0.325 & 0.377 & 0.500 & 0.652 & 0.640 & 0.653 & 0.624 & 0.691 & 0.720 \\ 0.289 & 0.329 & 0.348 & 0.500 & 0.666 & 0.659 & 0.656 & 0.661 & 0.771 \\ 0.298 & 0.266 & 0.360 & 0.334 & 0.500 & 0.631 & 0.709 & 0.671 & 0.686 \\ 0.256 & 0.298 & 0.347 & 0.341 & 0.369 & 0.500 & 0.668 & 0.738 & 0.669 \\ 0.298 & 0.347 & 0.376 & 0.344 & 0.291 & 0.332 & 0.500 & 0.734 & 0.617 \\ 0.341 & 0.353 & 0.309 & 0.339 & 0.329 & 0.262 & 0.266 & 0.500 & 0.646 \\ 0.291 & 0.419 & 0.280 & 0.229 & 0.314 & 0.331 & 0.383 & 0.354 & 0.500 \end{bmatrix}$$

TABLE 10. z summary table
True Model WST, $N = 100$

Method	VI distance _{MAP}			VI distance _{VI lb}			WAIC		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0	0.42	0.1	0	0.42	0.32	-13237.20 30.89	-13352.68 30.77	13657.35 31.35
WST model	0	0.42	0.1	0	0.42	0.32	-13273.59 31.16	-13293.26 30.37	-13673.15 31.34

In table (??), we report some summary statistics for the parameter z . As above, we have columns (a),(b) and (c) representing the case $K = 3, 5, 9$ respectively.

The first indicator is the $VI\text{distance}_{\text{MAP}}$ computed between the true partition z^* and the point estimate \hat{z}^{MAP} obtained with the maximum a posteriori estimate (MAP).

The second one is $VI\text{distance}_{\text{VI lb}}$ computed between the true partition z^* and the point estimate $\hat{z}^{\text{VI lb}}$ obtained with the partition attaining the VI lower bound.

The third one is the WAIC estimate, along with its standard error below.

TABLE 11. POMM hyperparameters summary table
True Model WST, $N = 100$

Fitted Model	$\hat{\theta}$			95% CI		
	(a)	(b)	(c)	(a)	(b)	(c)
σ	0.19	0.18	0.15	[0.01 0.79]	[0.04 0.65]	[0.07 0.27]
α	0.60	0.41	0.23	[0.24 0.9]	[0.11 0.75]	[0.1 0.4]

In table (??), we report the results for the hyperparameters of the *POMM* model. In the $\hat{\theta}$ column we report the estimates both for α and σ , while on the right we have their 95% Credible Interval. Given that the data were generated according the WST model, we do not have a ground truth to which these results should be compared. However, we can still try to make sense of these values by inspecting the properties of the induced P^{POMM} matrix resulting from the estimates.

The WST model has a prior over $P^{\text{WST}} \sim \text{Beta}(1, 1)$, then we may expect P^{POMM} itself to get as closer as needed to a uniform distribution by selecting the appropriate combination of (α, σ) .

Therefore, we simulate $n = 10000$ P^{POMM} matrices via the estimated parameters $\hat{\theta}$ in (??). Then we extract 1000 points from each level set to avoid sample biases, and we compare them with an equally-sized set simulated via the WST model, using the Kolmogorov-Smirnov test, where the null hypothesis is that two sets of points are sampled from the same distribution.

TABLE 12. Kolmogorov-Smirnov test
Data are generated via the estimated parameters

Method	p-value			% overlap between level sets		
	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0.26	0.05	0.00	82%	86%	89%

In table (??) we report first the Kolmogorov-Smirnov test p-values, and we may notice that for $K = 3, 5$ we do not reject with an $\alpha = 5\%$ that the P^{POMM} is compatible with being extracted from the WST model.

Instead with $K = 9$ we cannot say that P^{POMM} has collapsed to a uniform, but at the same time, if we look at the area of overlap between the densities of the level sets, we may notice that there is a significant amount of overlap between them, allowing the POMM model to effectively replicate the Unordered entries of the P^{WST} matrix.

In figure (??) we report the densities of the points generated via the estimated hyper-parameters. We may notice the difference for the case $K = 9$ with respect to the other two. In this case, we have very significant distributions for the POMM and the WST one.

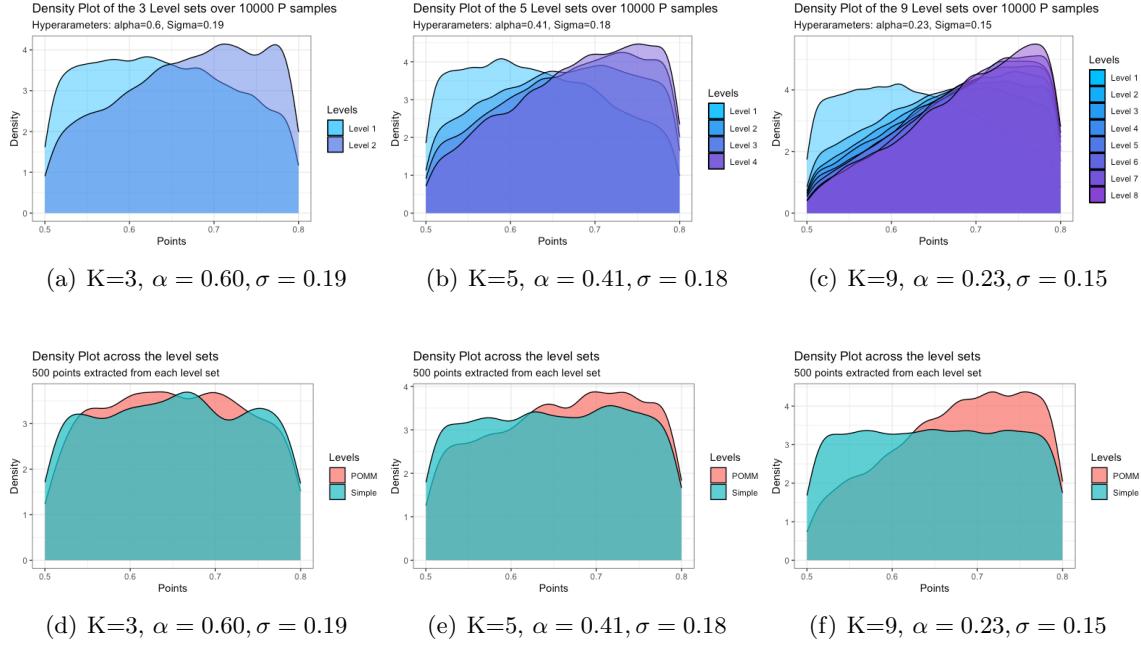


FIGURE 7. These are the densities of the entries of 10000 P matrices generated according to the parameters within brackets, that is, the parameters estimated according the POMM model on the data generated via the WST one. In figures (??), (??), (??) we have the densities of the level sets coloured differently. In figures (??), (??), (??) we put together 1000 points extracted from each level sets and we compute the density, so to have an overview of the joint distribution. We also compare the P 's entries simulated via the POMM and the WST model

10.0.1. *WST Model check.* In this subsection, we report some diagnostic checks for the algorithm of the Markov Chains, to assess convergence, quality of mixing, and the overall behaviour of the Metropolis-within-Gibbs algorithm.

TABLE 13. z diagnostic table
True Model WST, $N = 100$

Fitted Model	\overline{ESS}			\overline{ACF}_{30}			$\overline{\%accepted}$			$\overline{Gelman - Rubin}$		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	7558.60	7988.02	2880.01	0	0.01	0.08	0.02	0.04	0.15	1.29	1.03	1.31
WST	6544.85	6688.51	2438.17	0	0.03	0.11	0.01	0.56	0.15	1.04	1.29	1.26

In table (??) we report the diagnostics for the z parameter. Since the parameter is a label vector with 100 entries in this case, we compute the relevant statistics for each single entries an then we report the average.

- The first statistics is the Effective Sample Size (ESS) averaged over individuals $i = 1, \dots, N$, which denotes a fairly good sample size. The average is taken as follows:

$$\overline{ESS} = \frac{1}{n} \sum_{i=1}^N ESS_i$$

The same is applied also to the other diagnostic metrics.

- Then, we report the average autocorrelation, \overline{ACF}_{30} , computed with a lag of 30 iterations. This values are close to zero, meaning that there is very little correlation within the chain.
- In the third column, we report the average acceptance rate $\overline{\%accepted}$. These are significantly lower than the target acceptance rate that should be hit by the adaptive MCMC, which is 22%. However, the estimates are capable of correctly recovering the true partition. Combining the two facts, we may hypothesise that the simulation study has too "many" data, and proposing for a given individual i , who has already been assigned to the true block, a block different from the true one, leads to a drop in the likelihood which is too large, and as a consequence, we always reject the other labels.
- Finally, we compute the median Gelman-Rubin statistics for each entry, $\overline{Gelman - Rubin}$. Gelman and Rubin (1992) propose a general approach to monitoring convergence of MCMC output in which $m > 1$ parallel chains are run with starting values that are overdispersed relative to the posterior distribution. Convergence is diagnosed when the chains have 'forgotten' their initial values, and the output from all chains is indistinguishable. The Gelman-Rubin diagnostic is applied to a single variable

from the chain. It is based a comparison of within-chain and between-chain variances, and is similar to a classical analysis of variance. Values substantially above 1 indicate lack of convergence. If the chains have not converged,

TABLE 14. P diagnostic table
True Model WST, $N = 100$

Fitted Model	\overline{ESS}			$\overline{ACF_{30}}$			$\overline{\% \text{accepted}}$			$Gelman - Rubin$		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	2234.00	2101.1	1454.08	0.02	0.00	0.15	34.52	31.55	29.98	1.06	1.0	1.04
WST	2742.67	1446.9	1565.14	0.00	0.13	0.17	36.85	31.01	30.06	1.00	1.9	1.03

In table (??) we report the same diagnostics checks for the z parameter. The only difference is that here we do not average the diagnostics indicators over the individuals $i = 1, \dots, N$, but instead over the upper-triangular P indices: $\{i = 1, \dots, K - 1, \quad j = i + 1, \dots, K\}$.

Just as an example, the average ESS in this case is obtained as

$$\overline{ESS} = \frac{1}{(K \cdot (K - 1))/2} \sum_{i=1}^{N-1} \sum_{j=i+1}^N ESS_{i,j}$$

TABLE 15. POMM hyperparameters diagnostic table
True Model WST, $N = 100$

Fitted Model	ESS			ACF_{30}			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
σ	985	155	162	0.5	0.55	0.55	31.47	29.96	29.81	1.96	1.02	1.18
α	13	17	45	0.95	0.94	0.84	24.74	24.93	24.38	1.21	1.62	1.01

Finally, in table (??) we report again the same diagnostics, but since both α and σ are one-dimensional, we are presenting the diagnostics itself, without any average.

11. SIMULATION STUDY FROM THE POMM MODEL N=100

In this section we reverse the exercise performed in previous one. Before we were simulating from the WST model, now we are simulating from the POMM, with $K = 3, 5, 9$. The metrics, indices and summaries are the same as before, so we avoid replicating the same explanations.

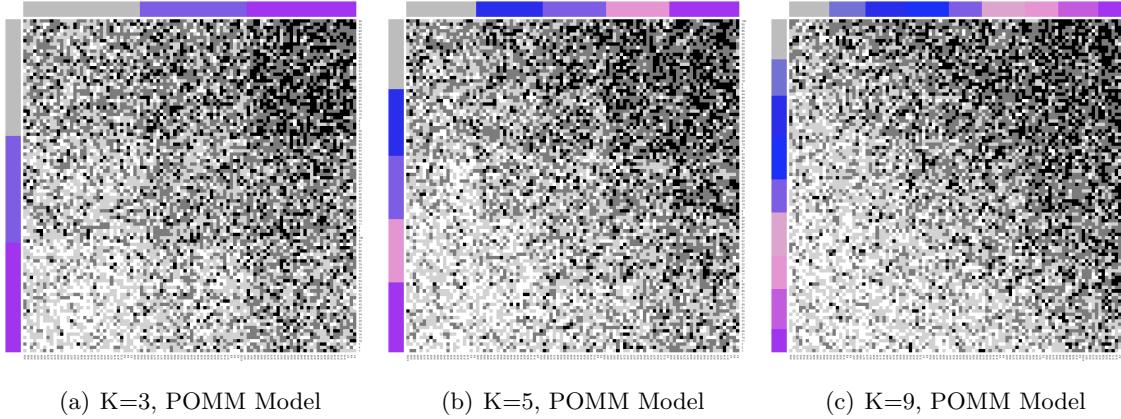


FIGURE 8. Adjacency Matrices simulated via the POMM Model

P summary table
True Model POMM, $K = 3, N = 100$

Fitted Model	\overline{MAE}			% within-95%-CI interval			CI interval length		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0.02	0.02	0.01	100%	90%	61.11%	0.14	0.08	0.02
WST model	0.00	0.04	0.02	100%	90%	100.00%	0.02	0.11	0.15

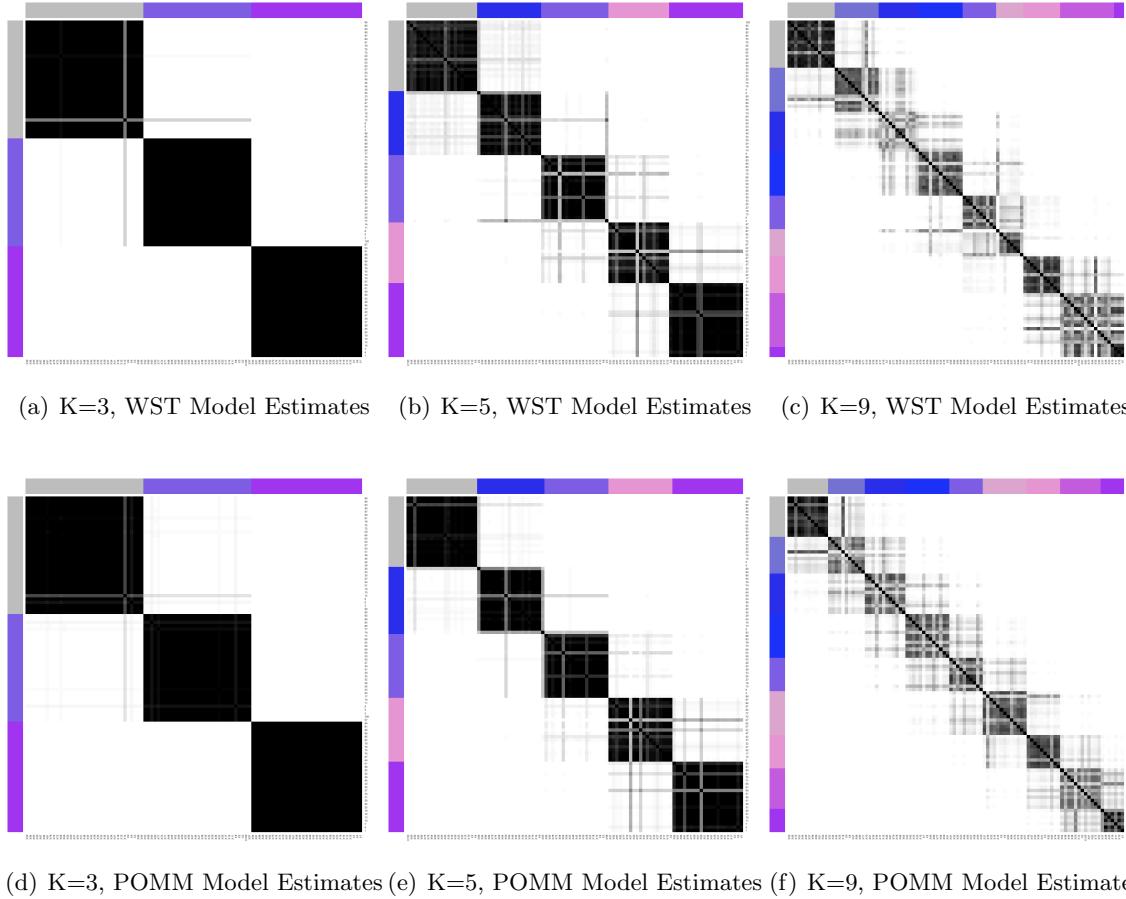


FIGURE 9. Co-Clustering Matrices obtained via the WST Model (above) and the POMM model (below).

$$P_{true}^{K=3} = \begin{bmatrix} 0.500 & 0.600 & 0.754 \\ 0.400 & 0.500 & 0.622 \\ 0.246 & 0.378 & 0.500 \end{bmatrix} \quad \hat{P}_{POMM}^{K=3} = \begin{bmatrix} 0.500 & 0.582 & 0.726 \\ 0.418 & 0.500 & 0.649 \\ 0.274 & 0.351 & 0.500 \end{bmatrix}$$

$$\hat{P}_{Simple}^{K=3} = \begin{bmatrix} 0.500 & 0.606 & 0.758 \\ 0.394 & 0.500 & 0.623 \\ 0.242 & 0.377 & 0.500 \end{bmatrix}$$

$$\begin{aligned}
P_{true}^{K=5} &= \begin{bmatrix} 0.500 & 0.569 & 0.679 & 0.752 & 0.781 \\ 0.431 & 0.500 & 0.576 & 0.698 & 0.741 \\ 0.321 & 0.424 & 0.500 & 0.562 & 0.674 \\ 0.248 & 0.302 & 0.438 & 0.500 & 0.571 \\ 0.219 & 0.259 & 0.326 & 0.429 & 0.500 \end{bmatrix} & \hat{P}_{POMM}^{K=5} &= \begin{bmatrix} 0.500 & 0.575 & 0.694 & 0.723 & 0.775 \\ 0.425 & 0.500 & 0.604 & 0.655 & 0.736 \\ 0.306 & 0.396 & 0.500 & 0.555 & 0.655 \\ 0.277 & 0.345 & 0.445 & 0.500 & 0.597 \\ 0.225 & 0.264 & 0.345 & 0.403 & 0.500 \end{bmatrix} \\
&& \hat{P}_{Simple}^{K=5} &= \begin{bmatrix} 0.500 & 0.557 & 0.681 & 0.703 & 0.761 \\ 0.443 & 0.500 & 0.659 & 0.641 & 0.744 \\ 0.319 & 0.341 & 0.500 & 0.532 & 0.625 \\ 0.297 & 0.359 & 0.468 & 0.500 & 0.622 \\ 0.239 & 0.256 & 0.375 & 0.378 & 0.500 \end{bmatrix} \\
P_{true}^{K=9} &= \begin{bmatrix} 0.500 & 0.547 & 0.626 & 0.682 & 0.699 & 0.726 & 0.766 & 0.775 & 0.778 \\ 0.453 & 0.500 & 0.546 & 0.624 & 0.679 & 0.702 & 0.729 & 0.750 & 0.765 \\ 0.374 & 0.454 & 0.500 & 0.571 & 0.633 & 0.647 & 0.705 & 0.720 & 0.738 \\ 0.318 & 0.376 & 0.429 & 0.500 & 0.551 & 0.618 & 0.660 & 0.692 & 0.708 \\ 0.301 & 0.321 & 0.367 & 0.449 & 0.500 & 0.561 & 0.630 & 0.655 & 0.710 \\ 0.274 & 0.298 & 0.353 & 0.382 & 0.439 & 0.500 & 0.557 & 0.625 & 0.676 \\ 0.234 & 0.271 & 0.295 & 0.340 & 0.370 & 0.443 & 0.500 & 0.562 & 0.636 \\ 0.225 & 0.250 & 0.280 & 0.308 & 0.345 & 0.375 & 0.438 & 0.500 & 0.560 \\ 0.222 & 0.235 & 0.262 & 0.292 & 0.290 & 0.324 & 0.364 & 0.440 & 0.500 \end{bmatrix} \\
&& \hat{P}_{POMM}^{K=9} &= \begin{bmatrix} 0.500 & 0.559 & 0.637 & 0.674 & 0.704 & 0.729 & 0.754 & 0.773 & 0.791 \\ 0.441 & 0.500 & 0.558 & 0.637 & 0.675 & 0.704 & 0.730 & 0.751 & 0.772 \\ 0.363 & 0.442 & 0.500 & 0.559 & 0.637 & 0.674 & 0.704 & 0.730 & 0.752 \\ 0.326 & 0.363 & 0.441 & 0.500 & 0.557 & 0.636 & 0.675 & 0.704 & 0.729 \\ 0.296 & 0.325 & 0.363 & 0.443 & 0.500 & 0.557 & 0.638 & 0.676 & 0.705 \\ 0.271 & 0.296 & 0.326 & 0.364 & 0.443 & 0.500 & 0.557 & 0.637 & 0.675 \\ 0.246 & 0.270 & 0.296 & 0.325 & 0.362 & 0.443 & 0.500 & 0.558 & 0.638 \\ 0.227 & 0.249 & 0.270 & 0.296 & 0.324 & 0.363 & 0.442 & 0.500 & 0.558 \\ 0.209 & 0.228 & 0.248 & 0.271 & 0.295 & 0.325 & 0.362 & 0.442 & 0.500 \end{bmatrix} \\
&& \hat{P}_{Simple}^{K=9} &= \begin{bmatrix} 0.500 & 0.547 & 0.626 & 0.682 & 0.699 & 0.726 & 0.766 & 0.775 & 0.778 \\ 0.453 & 0.500 & 0.546 & 0.624 & 0.679 & 0.702 & 0.729 & 0.750 & 0.765 \\ 0.374 & 0.454 & 0.500 & 0.571 & 0.633 & 0.647 & 0.705 & 0.720 & 0.738 \\ 0.318 & 0.376 & 0.429 & 0.500 & 0.551 & 0.618 & 0.660 & 0.692 & 0.708 \\ 0.301 & 0.321 & 0.367 & 0.449 & 0.500 & 0.561 & 0.630 & 0.655 & 0.710 \\ 0.274 & 0.298 & 0.353 & 0.382 & 0.439 & 0.500 & 0.557 & 0.625 & 0.676 \\ 0.234 & 0.271 & 0.295 & 0.340 & 0.370 & 0.443 & 0.500 & 0.562 & 0.636 \\ 0.225 & 0.250 & 0.280 & 0.308 & 0.345 & 0.375 & 0.438 & 0.500 & 0.560 \\ 0.222 & 0.235 & 0.262 & 0.292 & 0.290 & 0.324 & 0.364 & 0.440 & 0.500 \end{bmatrix}
\end{aligned}$$

11.1. POMM model check.

z summary table
True Model POMM, $N = 100$

Method	VI distanceMAP			VI distanceVI lb			WAIC		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0.13	0.53	1.9	0.13	0.42	1.73	-13361.79 30.98	-13510.00 31.45	-13645.28 31.03
Simple model	0.13	0.31	2.0	0.13	0.48	1.71	-13409.44 30.74	-13496.19 31.71	-13659.10 30.71

TABLE 16. POMM Hyperparameters summary table
True Model POMM, $N = 100$

Method	$\hat{\theta}$			95% CI interval			True value		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
σ	0.1487	0.084	0.0196	[0.0033 0.7575]	[4e-04 0.5602]	[0.0012 0.0625]	0.01	0.01	0.01
α	0.4237	0.5265	0.491	[0.1324 0.6215]	[0.439 0.8733]	[0.4258 0.5934]	0.5	0.5	0.5

z diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	3418.94	10354.90	4800.38	0.11	0.02	0.09	0.01	0.33	4.50	1.29	1.08	1.00
Simple	1778.56	10191.58	1869.04	0.00	0.01	0.36	0.01	0.50	2.64	1.02	1.17	1.08

P diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	626.67	672.7	308.22	0.29	0.22	0.34	20.89	29.75	29.55	1.36	1.34	1.11
Simple	2798.00	1659.7	121.94	0.00	0.02	0.53	36.70	30.99	29.78	1.00	1.07	1.03

POMM hyperparameters diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
σ	981	480	13	0.4	0.58	0.94	15.78	11.38	1.72	3.55	2.15	1.96
α	59	85	29	0.8	0.71	0.88	16.32	15.41	6.33	1.26	1.35	1.25

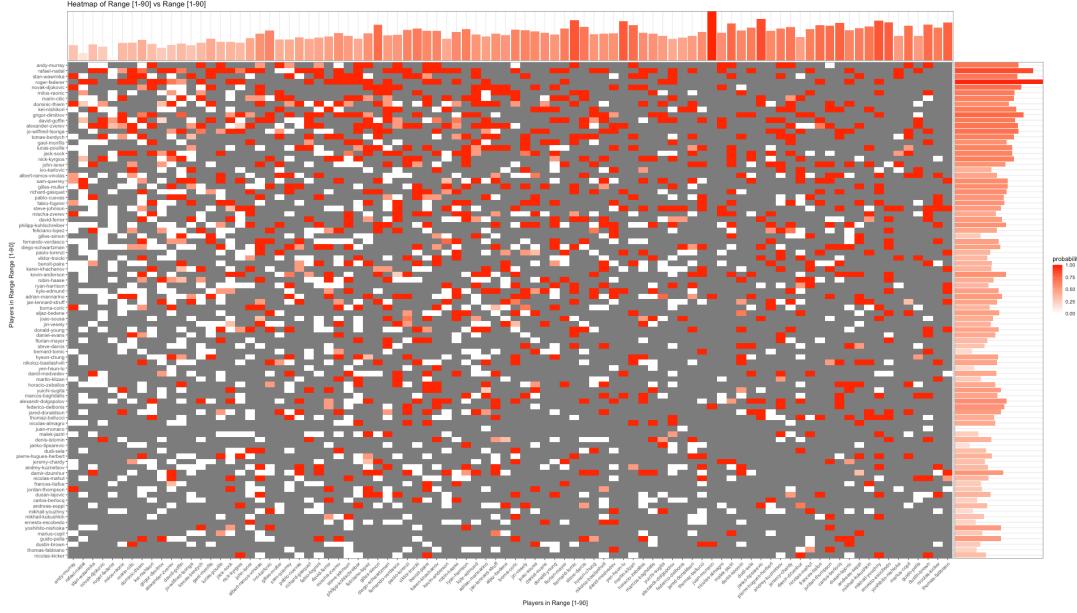


FIGURE 10. Adjacency matrix of the top 90 ranked players. Players on the x -axis are ordered in ascending rank order (from 1st to 90th), while one the y -axis they are sorted in descending order (from 90th to 1st). Each entry of the matrix above is coloured according to a scale of colours from white to red, which gets toward the red as the ratio $\hat{p}_{ij} = \frac{\text{victories}_{ij}}{\text{matches}_{ij}}$, where i is the row index and j is the column index, grows. On the sides we have the marginal probabilities. On the x -axis, we have bars that are proportional to $\bar{p}_{ij^*} = \frac{\sum_{i=1}^{90} \text{victories}_{ij^*}}{\sum_{i=1}^{90} \text{matches}_{ij^*}}$. On the y -axis, we have bars that are proportional to $\bar{p}_{i^*j} = \frac{\sum_{j=1}^{90} \text{victories}_{i^*j}}{\sum_{j=1}^{90} \text{matches}_{i^*j}}$.

12. EXPLORATORY ANALYSIS OF THE TENNIS DATA

12.1. Overview of data. In figure (??) we have the adjacency matrix representing the interactions between the top 90 ranked players. The x -axis lists players in ascending rank order (from 1st to 90th), while the y -axis lists players in descending order (from 90th to 1st). Each cell in the matrix is shaded according to a color scale from white to red, indicating the ratio $\hat{p}_{ij} = \frac{\text{victories}_{ij}}{\text{matches}_{ij}}$. This ratio, where i is the row index and j is the column index, becomes redder as the player i has more victories relative to matches against player j .

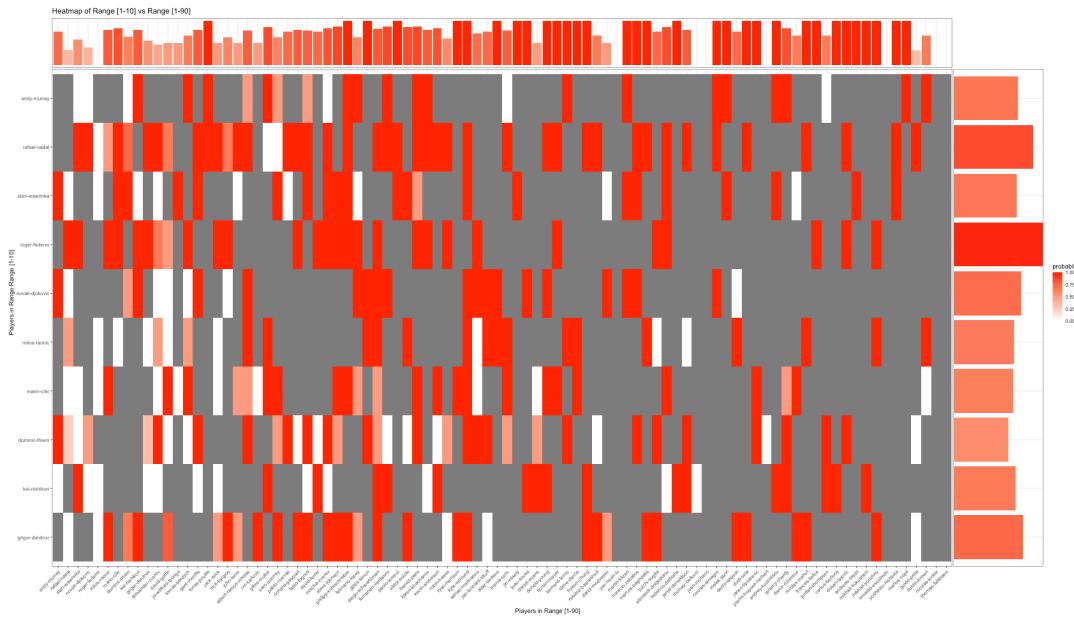
The side margins display marginal probabilities. On the x -axis, we have bars proportional to $\bar{p}_{ij^*} = \frac{\sum_{i=1}^{90} \text{victories}_{ij^*}}{\sum_{i=1}^{90} \text{matches}_{ij^*}}$, representing the average probability of player i winning against player j^* across all players.

On the y -axis, we have bars proportional to $\bar{p}_{i^*j} = \frac{\sum_{j=1}^{90} \text{victories}_{i^*j}}{\sum_{j=1}^{90} \text{matches}_{i^*j}}$, representing the average probability of player i^* winning against player j across all players.

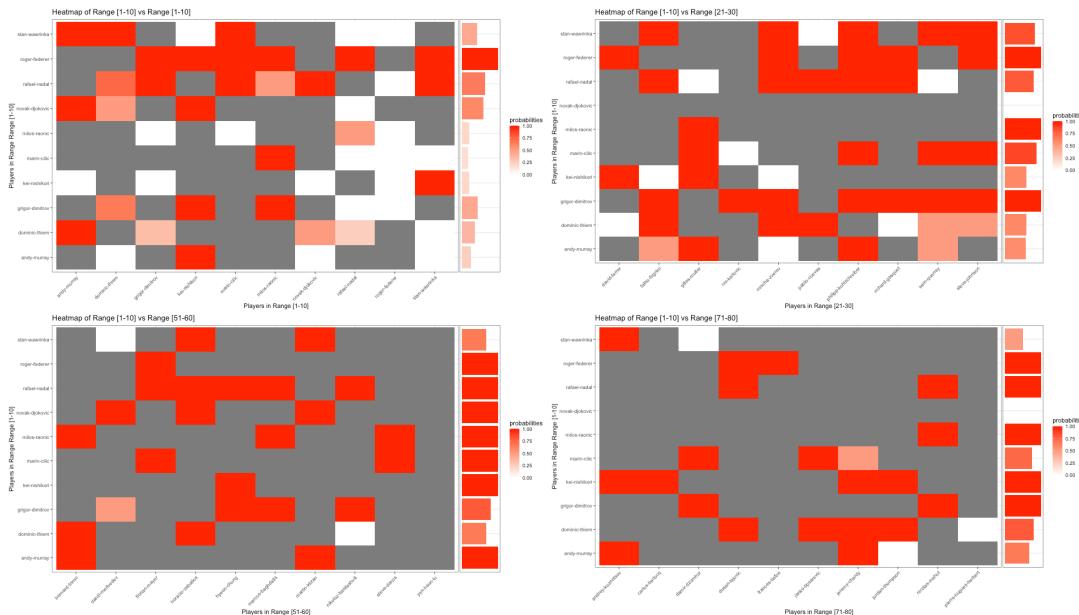
12.2. Focus on the top 10 Ranked Players. Here we take a closer look at the players ranked from 1 to 10.

Range [1-10].csv

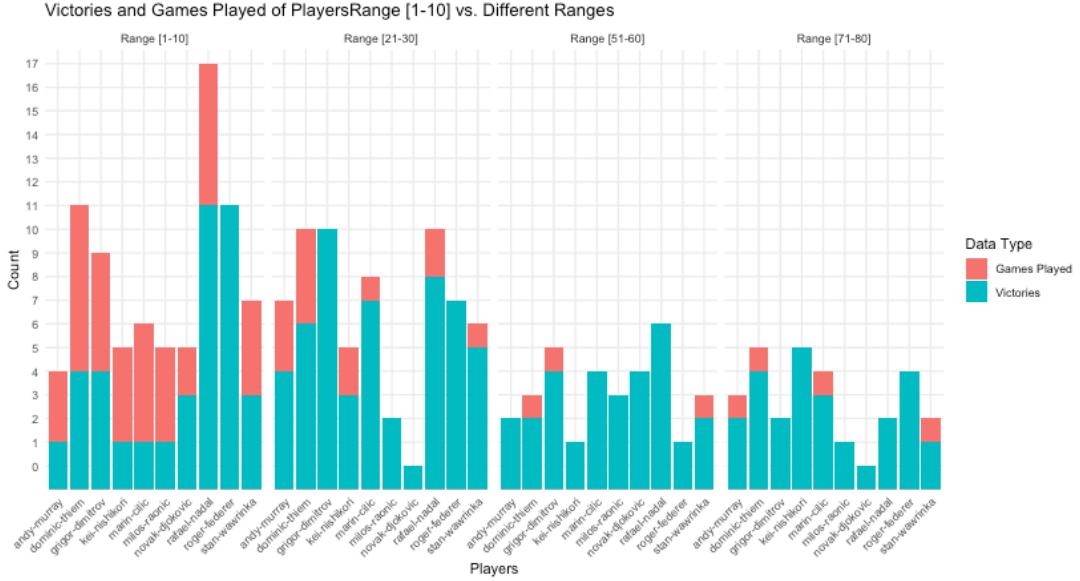
Range	Rank players	Mean n_{Matches}	Mean $n_{\text{Victories}}$	Median $\frac{n_{\text{Victories}}}{n_{\text{Matches}}}$	IQR Median	sd Median
Rank [1-10]		6.5	3.0	0.3961039	0.3486111	0.2609968
Rank [21-30]		7.0	5.5	0.8333333	0.4000000	0.1796782
Rank [51-60]		3.0	2.5	1.0000000	0.1500000	0.1441878
Rank [71-80]		2.5	2.0	1.0000000	0.2500000	0.1872890



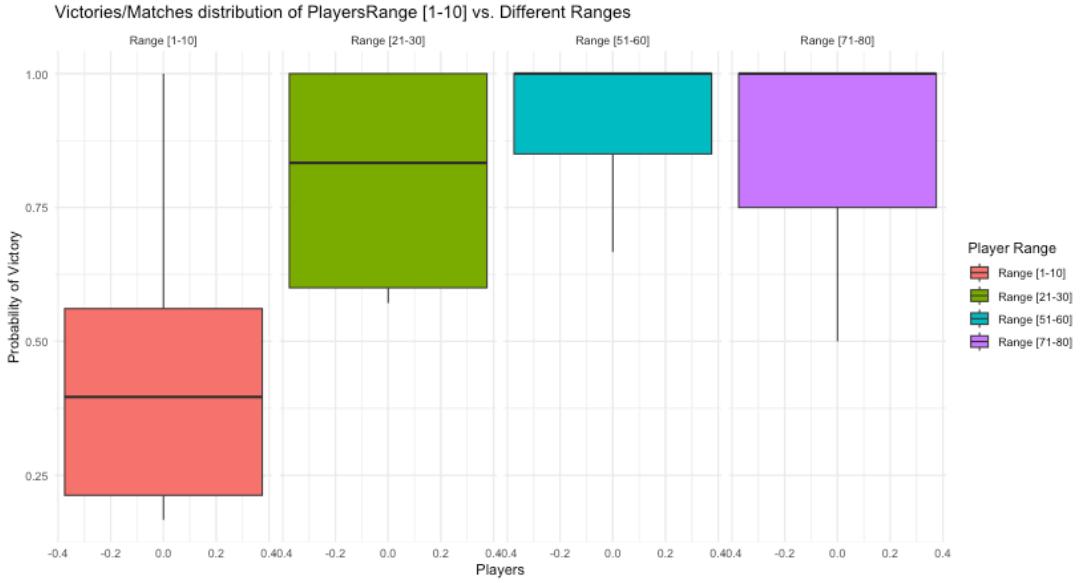
(a) Adjacency of 1-10th ranked players versus the other 1-90th players.



(b) Adjacency of 1-10th ranked players versus 1-10th players (upper-left quadrant), 21-30th players (upper-right quadrant), 51-60th players (bottom-left quadrant), 71-80th players (bottom-right quadrant)



(c) Bar plot comparing the performance of players ranked 1-10 with other player groups. The x-axis displays the 1st-10th ranked players repeated four times. Each panel provides data on the number of victories and total number of games played. These statistics are compared between the 1-10th ranked players and players ranked 21-30, 51-60, and 71-80. In each panel, the red bar represents the total number of games played by each player in the 1-10 position against all players in one of the specific rank blocks mentioned above. The green bar corresponds to the number of victories achieved by the same player in these matchups.

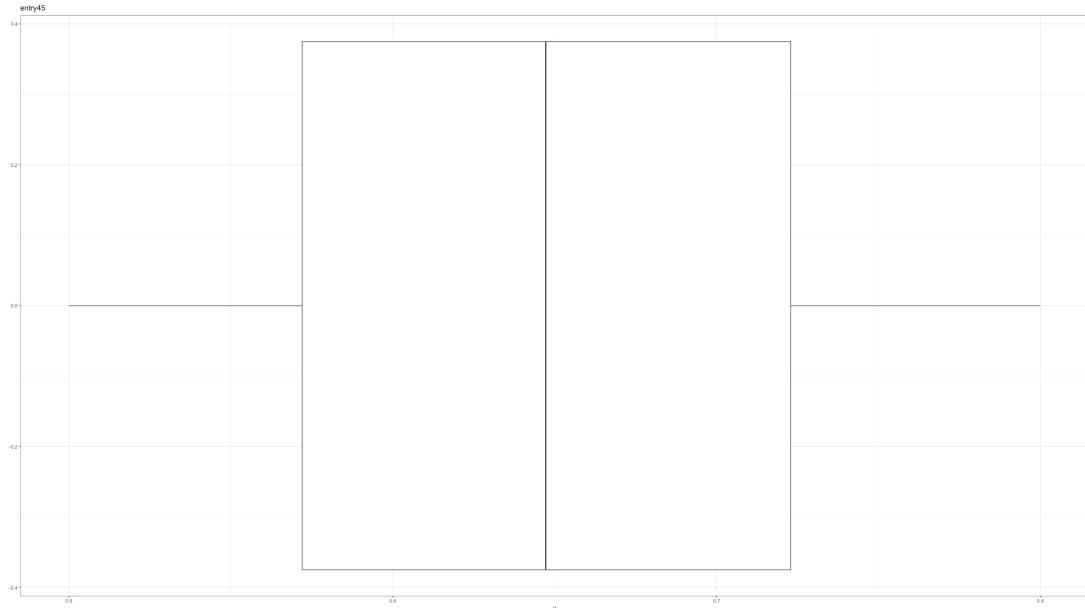


(d) This box plot illustrates the distribution of the quantity $\frac{\text{victories}}{\text{games}}$ for players ranked 1-10 playing against other players ranked 21-30, 51-60, and 71-80, arranged from left to right. The vertical axis represents the ratio of victories to games of the players in the group 1-10, providing insights into the performance of these players across various rank ranges. The boxes indicate the interquartile range (IQR) of the data, while the whiskers extend to the minimum and maximum values within 1.5 times the IQR. Outliers, if present, are shown as individual data points beyond this range.

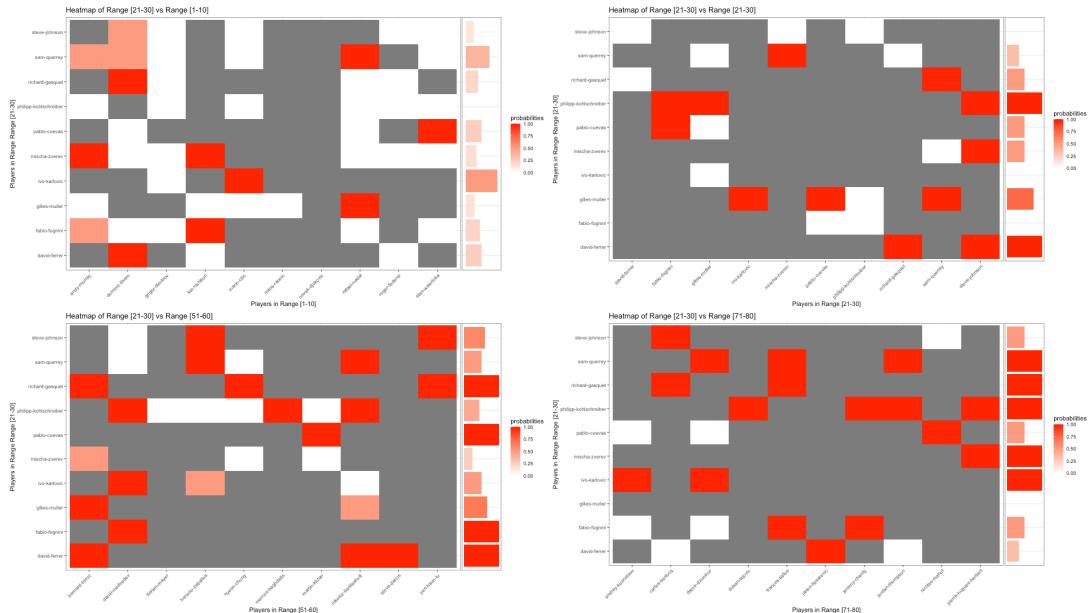
12.3. Focus on the 21st-30th Ranked Players. Here we take a closer look at the players ranked from the 21st position to the 30th.

Range [1-10].csv

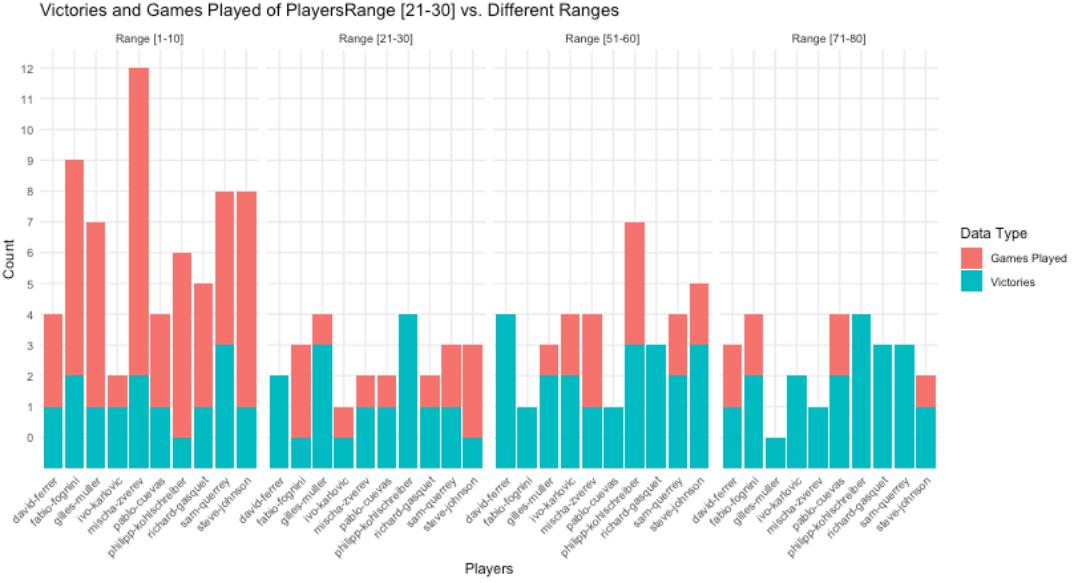
Range Rank players	Mean n_{Matches}	Mean $n_{\text{Victories}}$	Median $\frac{n_{\text{Victories}}}{n_{\text{Matches}}}$	IQR Median	sd Median
Range [1-10]	6.5	1	0.21	0.10	0.14
Range [21-30]	2.5	1	0.50	0.60	0.38
Range [51-60]	4.0	2	0.63	0.50	0.28
Range [71-80]	3.0	2	1.00	0.50	0.29



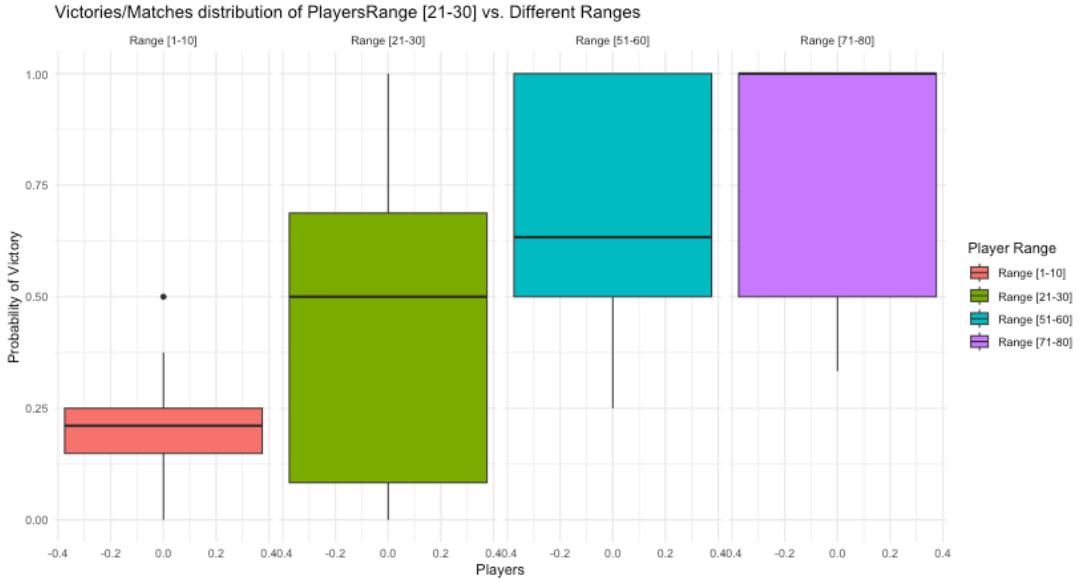
(e) Adjacency of 21-30th ranked players versus the other 1-90th players.



(f) Adjacency of 21-30th ranked players versus 1st-10th players (upper-left quadrant), 21-30th players (upper-right quadrant), 51-60th players (bottom-left quadrant), 71-80th players (bottom-right quadrant)



(g) Bar plot comparing the performance of players ranked 21-30 with other player groups. The x-axis displays the 21st-30th ranked players repeated four times. Each panel provides data on the number of victories and total number of games played. These statistics are compared between the 21st-30th ranked players and players ranked 1-10, 21-30, 51-60, and 71-80. In each panel, the red bar represents the total number of games played by each player in the 21-30 position against all players in one of the specific rank blocks mentioned above. The green bar corresponds to the number of victories achieved by the same player in these matchups.

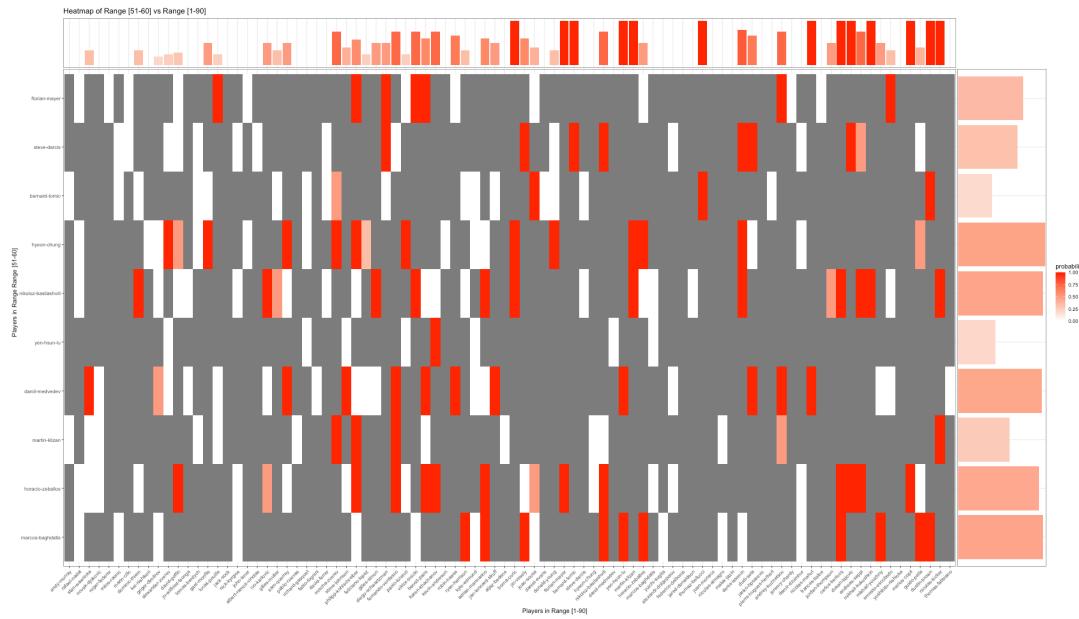


(h) This box plot illustrates the distribution of the quantity $\frac{\text{victories}}{\text{games}}$ for players ranked 21-30 playing against other players ranked 1-10, 21-30, 51-60, and 71-80, arranged from left to right. The vertical axis represents the ratio of victories to games of the players in the group 21-30, providing insights into the performance of these players across various rank ranges. The boxes indicate the interquartile range (IQR) of the data, while the whiskers extend to the minimum and maximum values within 1.5 times the IQR. Outliers, if present, are shown as individual data points beyond this range.

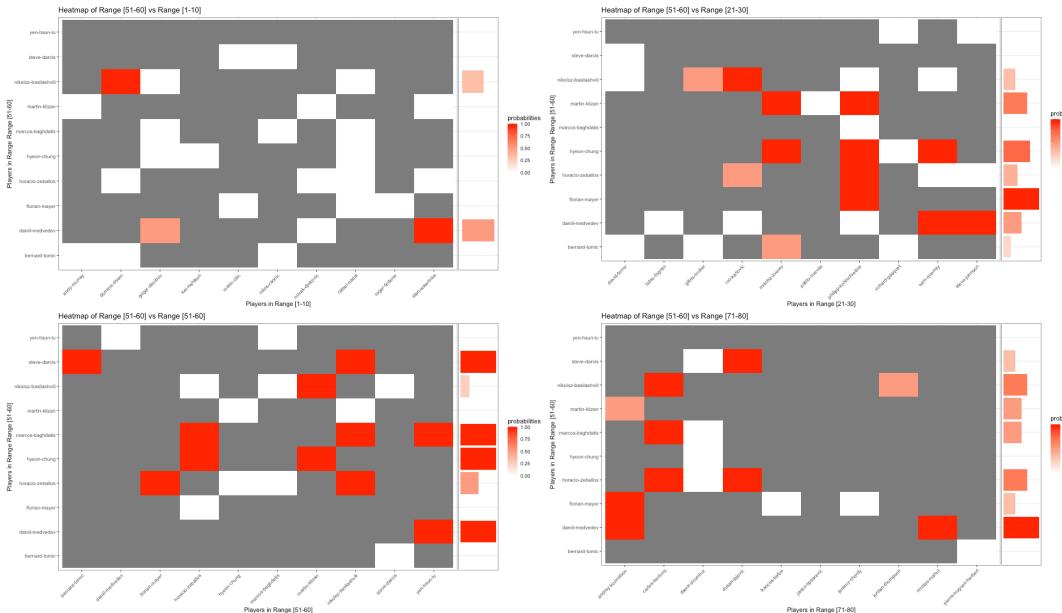
12.4. Focus on the 51st-60th Ranked Players. Here we take a closer look at the players ranked from the 51th to the 60th position.

Range [51-60].csv

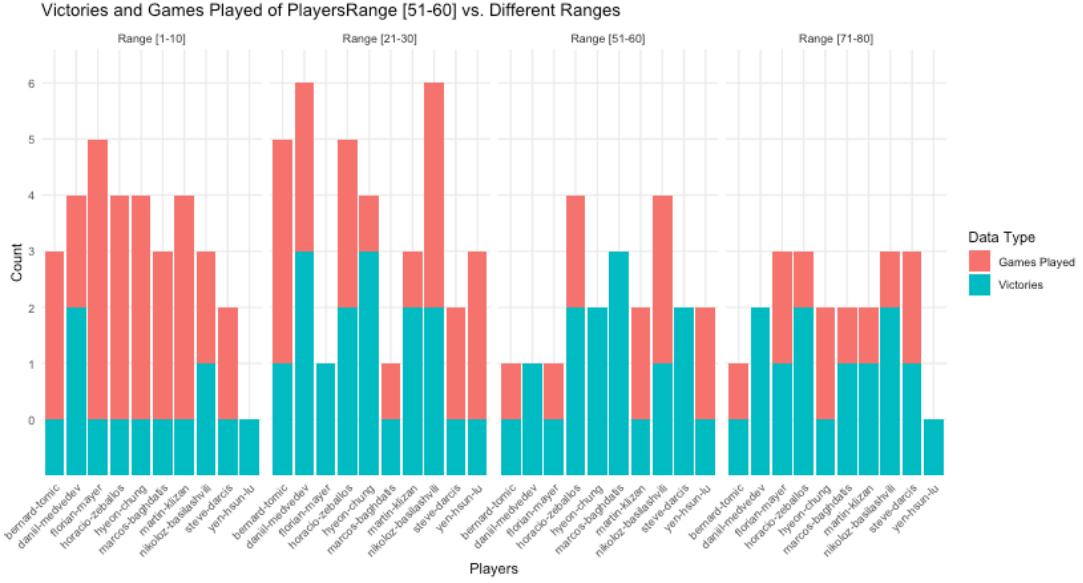
Range	Rank players	Mean n_{Matches}	Mean $n_{\text{Victories}}$	Median $\frac{n_{\text{Victories}}}{n_{\text{Matches}}}$	IQR Median	sd	Median
Range [1-10]		3.5	0.0	0.00	0.00	0.19	
Range [21-30]		3.5	1.5	0.37	0.58	0.35	
Range [51-60]		2.0	1.0	0.38	1.00	0.48	
Range [71-80]		2.0	1.0	0.50	0.33	0.32	



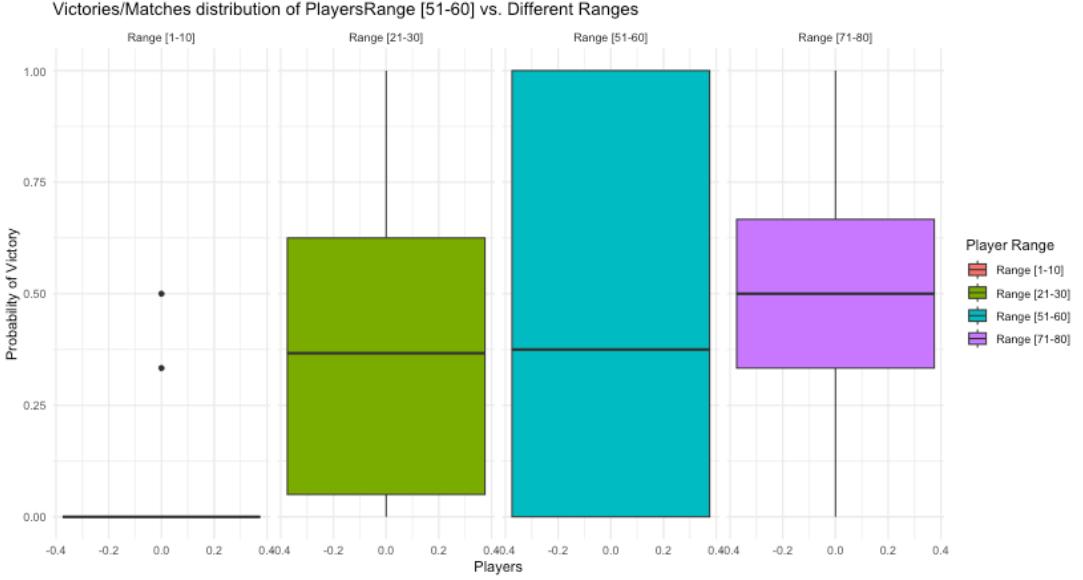
(i) Adjacency of 51-60th ranked players versus the other 1-90th players.



(j) Adjacency of 51-60th ranked players versus 1-10th players (upper-left quadrant), 21-30th players (upper-right quadrant), 51-60th players (bottom-left quadrant), 71-80th players (bottom-right quadrant)



(k) Bar plot comparing the performance of players ranked 51-60 with other player groups. The x-axis displays the 51st-60th ranked players repeated four times. Each panel provides data on the number of victories and total number of games played. These statistics are compared between the 51st-60th ranked players and players ranked 1-10, 21-30, 51-60, and 71-80. In each panel, the red bar represents the total number of games played by each player in the 51-60 position against all players in one of the specific rank blocks mentioned above. The green bar corresponds to the number of victories achieved by the same player in these matchups.

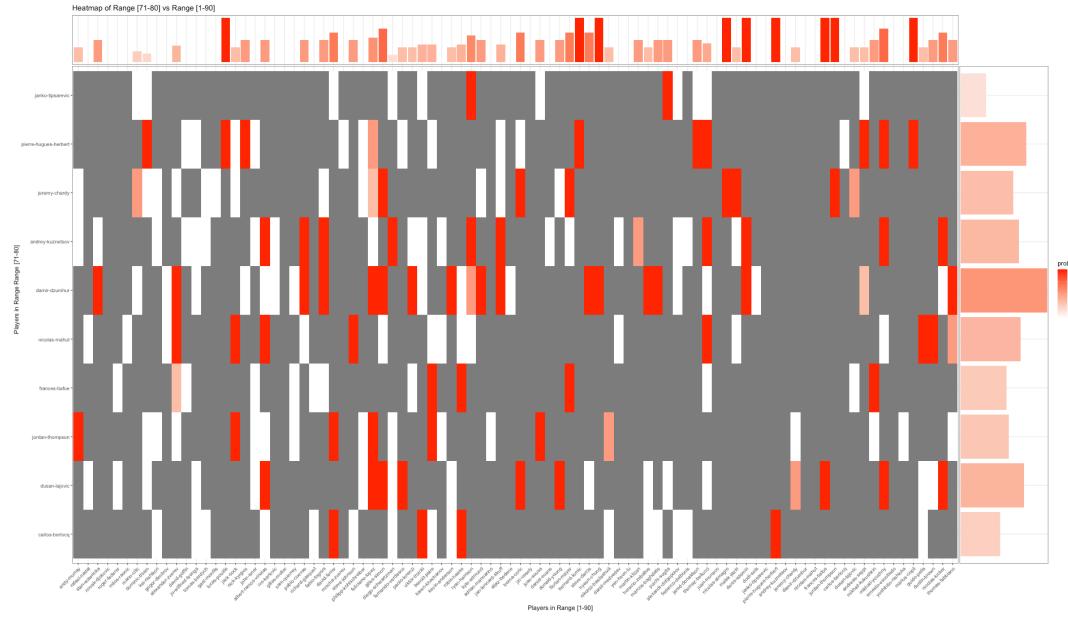


(l) This box plot illustrates the distribution of the quantity $\frac{\text{victories}}{\text{games}}$ for players ranked 51-60 playing against other players ranked 1-10, 21-30, 51-60, and 71-80, arranged from left to right. The vertical axis represents the ratio of victories to games of the players in the group 51-60, providing insights into the performance of these players across various rank ranges. The boxes indicate the interquartile range (IQR) of the data, while the whiskers extend to the minimum and maximum values within 1.5 times the IQR. Outliers, if present, are shown as individual data points beyond this range.

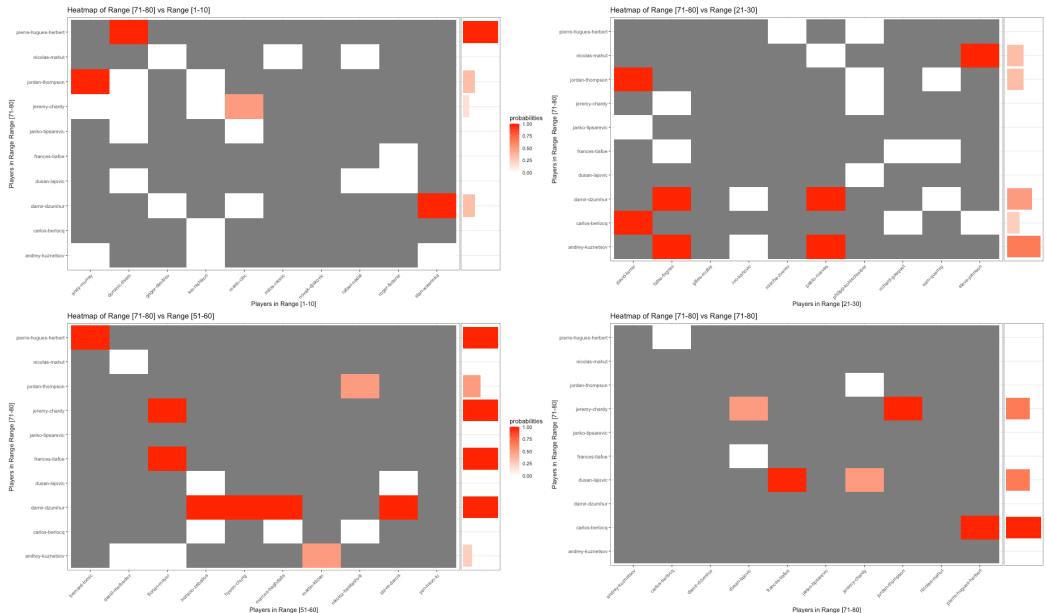
12.5. Focus on the 71st-80th Ranked Players. Here we take a closer look at the players ranked from 1 to 10.

Performance of players ranked within the interval [71-80] vs other rank groups

Range	Rank players	Mean n_{Matches}	Mean $n_{\text{Victories}}$	Median $\frac{n_{\text{Victories}}}{n_{\text{Matches}}}$	IQR	Median	sd	Median
Range [1-10]		3.0	0.0	0.00		0.29		0.32
Range [21-30]		3.0	0.5	0.13		0.33		0.25
Range [51-60]		1.5	1.0	0.50		1.00		0.48
Range [71-80]		1.0	0.0	0.33		0.67		0.44

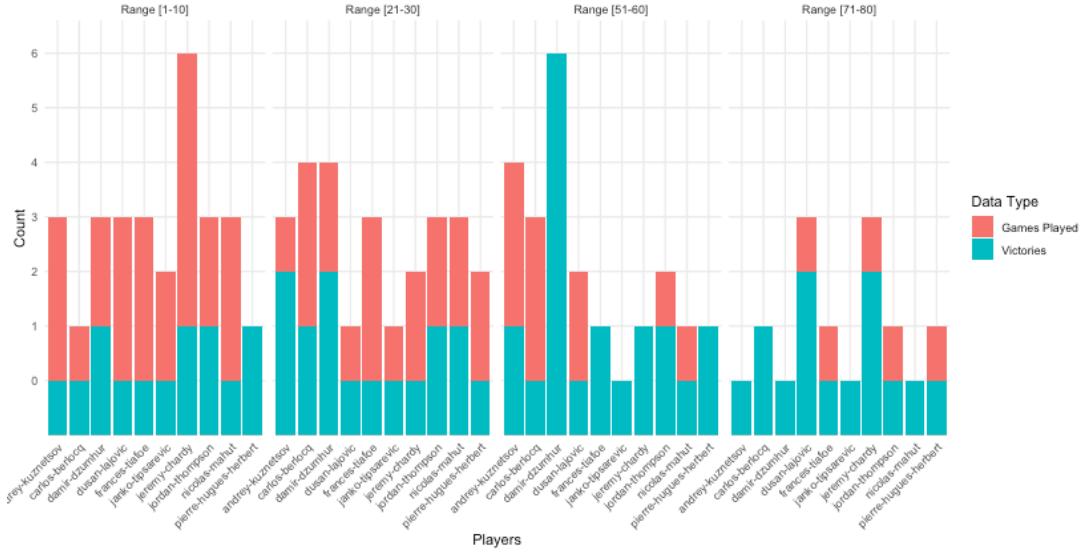


(m) Adjacency of 71-80th ranked players versus the other 1-90th players.



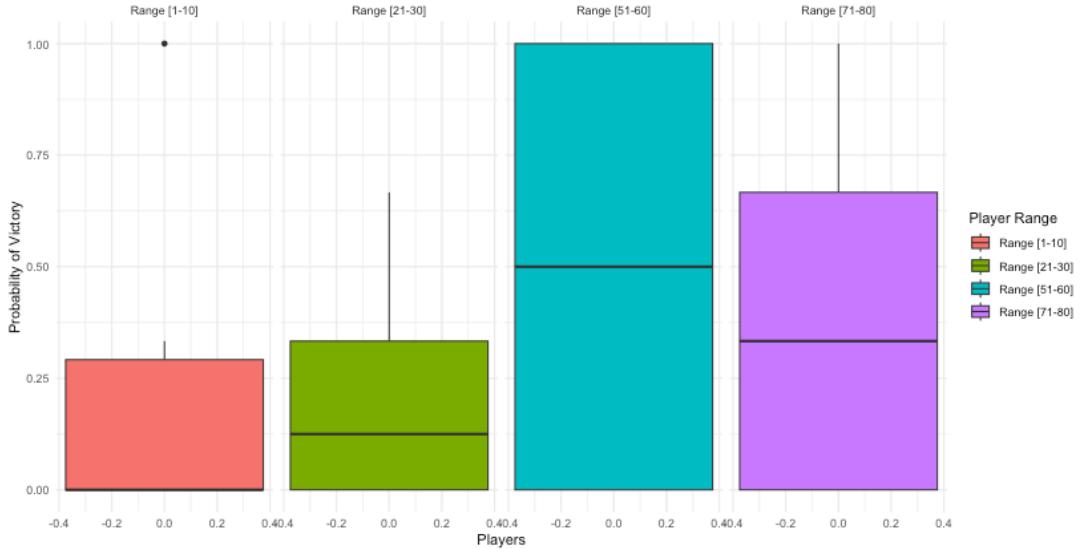
(n) Adjacency of 71-80th ranked players versus 1-10th players (upper-left quadrant), 21-30th players (upper-right quadrant), 51-60th players (bottom-left quadrant), 71-80th players (bottom-right quadrant)

Victories and Games Played of PlayersRange [71-80] vs. Different Ranges



(o) Bar plot comparing the performance of players ranked 71-80 with other player groups. The x-axis displays the 1st-10th ranked players repeated four times. Each panel provides data on the number of victories and total number of games played. These statistics are compared between the 1-10th ranked players and players ranked 21-30, 51-60, and 71-80. In each panel, the red bar represents the total number of games played by each player in the 71-80 position against all players in one of the specific rank blocks mentioned above. The green bar corresponds to the number of victories achieved by the same player in these matchups.

Victories/Matches distribution of PlayersRange [71-80] vs. Different Ranges



(p) This box plot illustrates the distribution of the quantity $\frac{\text{victories}}{\text{games}}$ for players ranked 71-80 playing against other players ranked 21-30, 51-60, and 71-80, arranged from left to right. The vertical axis represents the ratio of victories to games of the players in the group 71-80, providing insights into the performance of these players across various rank ranges. The boxes indicate the interquartile range (IQR) of the data, while the whiskers extend to the minimum and maximum values within 1.5 times the IQR. Outliers, if present, are shown as individual data points beyond this range.

13. APPLICATION TO TENNIS DATA

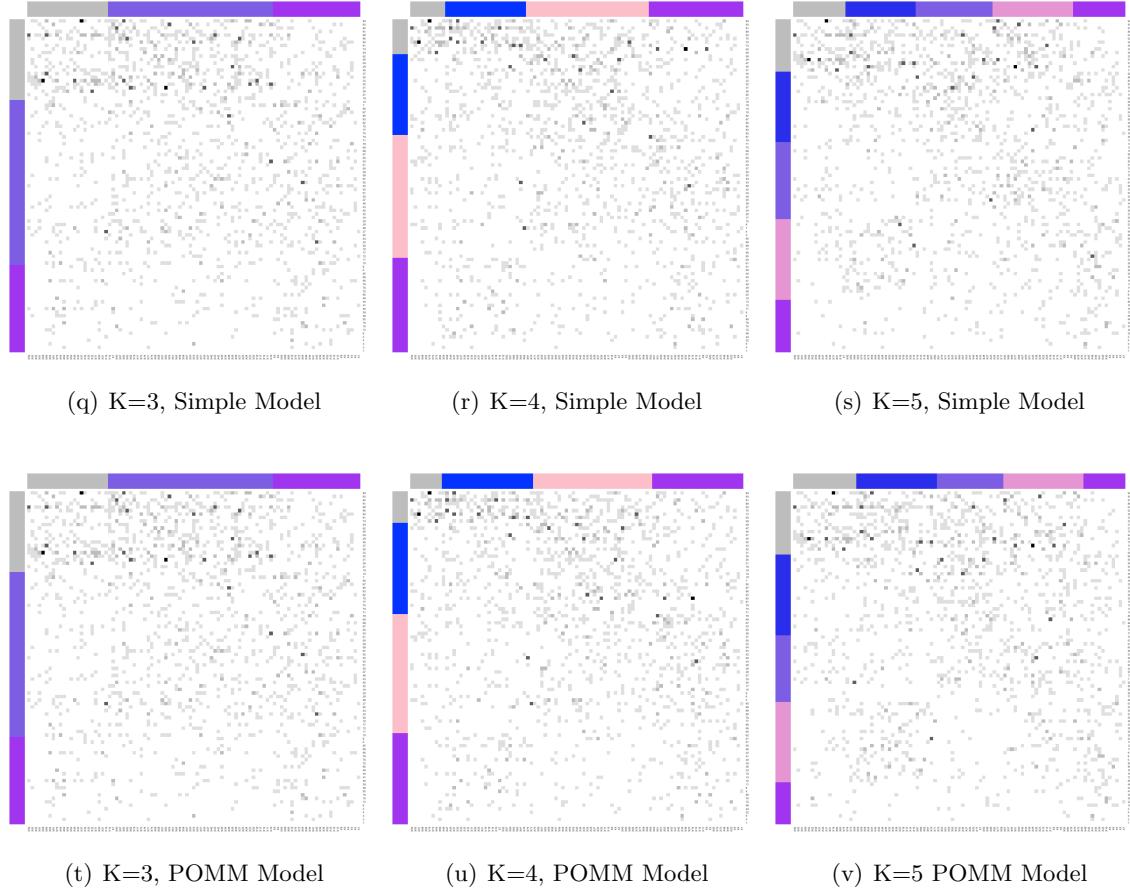


FIGURE 11. Adjacency Matrices simulated via the POMM Model

$$\hat{P}^{POMM} = \begin{bmatrix} 0.500 & 0.779 & 0.648 \\ 0.221 & 0.500 & 0.765 \\ 0.352 & 0.235 & 0.500 \end{bmatrix} \quad \hat{P}^{Simple} = \begin{bmatrix} 0.500 & 0.779 & 0.648 \\ 0.221 & 0.500 & 0.766 \\ 0.352 & 0.234 & 0.500 \end{bmatrix}$$

$$\hat{P}^{POMM} = \begin{bmatrix} 0.500 & 0.786 & 0.742 & 0.764 \\ 0.214 & 0.500 & 0.775 & 0.532 \\ 0.258 & 0.225 & 0.500 & 0.776 \\ 0.236 & 0.468 & 0.224 & 0.500 \end{bmatrix} \quad \hat{P}^{Simple} = \begin{bmatrix} 0.500 & 0.787 & 0.742 & 0.763 \\ 0.213 & 0.500 & 0.775 & 0.532 \\ 0.258 & 0.225 & 0.500 & 0.775 \\ 0.237 & 0.468 & 0.225 & 0.500 \end{bmatrix}$$

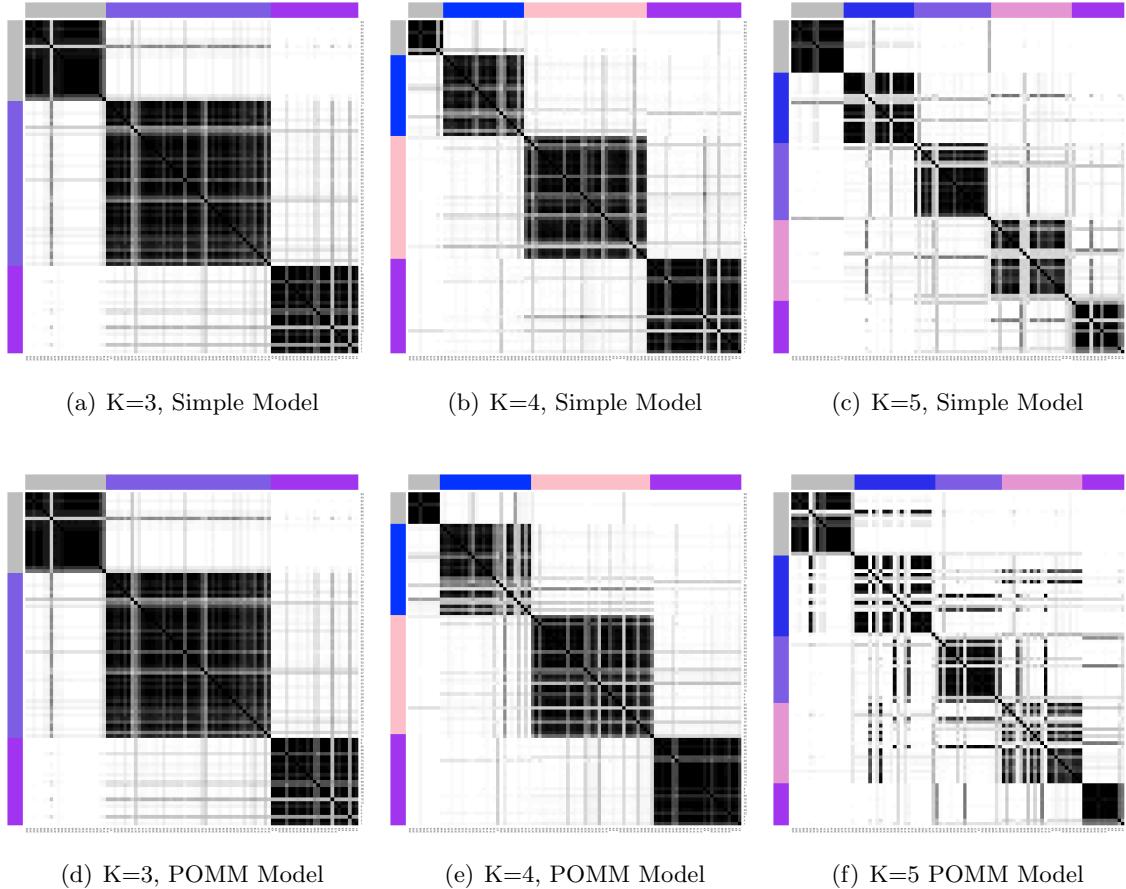


FIGURE 12. Adjacency Matrices simulated via the POMM Model

$$\hat{P}^{POMM} = \begin{bmatrix} 0.500 & 0.778 & 0.745 & 0.785 & 0.716 \\ 0.222 & 0.500 & 0.792 & 0.512 & 0.768 \\ 0.255 & 0.208 & 0.500 & 0.747 & 0.652 \\ 0.215 & 0.488 & 0.253 & 0.500 & 0.776 \\ 0.284 & 0.232 & 0.348 & 0.224 & 0.500 \end{bmatrix} \quad \hat{P}^{Simple} = \begin{bmatrix} 0.500 & 0.779 & 0.745 & 0.785 & 0.714 \\ 0.221 & 0.500 & 0.792 & 0.512 & 0.768 \\ 0.255 & 0.208 & 0.500 & 0.748 & 0.652 \\ 0.215 & 0.488 & 0.252 & 0.500 & 0.777 \\ 0.286 & 0.232 & 0.348 & 0.223 & 0.500 \end{bmatrix}$$

13.1. POMM model check.

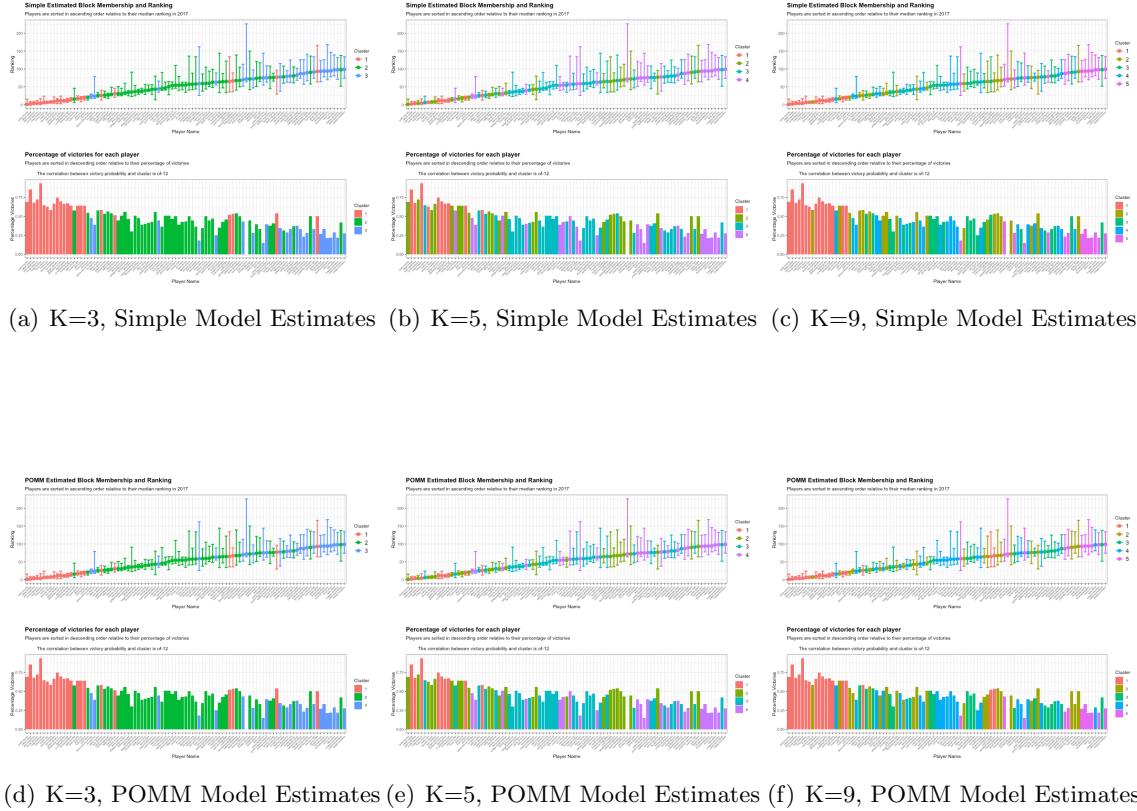


FIGURE 13. Co-Clustering Matrices obtained via the Simple Model (above) and the POMM model (below).

z summary table
True Model POMM, $N = 100$

Method	WAIC		
	(a)	(b)	(c)
POMM model	-5410.20 24.85	-5536.49 24.78	-5637.33 25.51
Simple model	-5411.18 24.87	-5535.89 24.76	-5637.28 25.48

POMM Hyperparameters summary table
True Model POMM, $N = 100$

Method	$\hat{\theta}$			95% CI interval		
	(a)	(b)	(c)	(a)	(b)	(c)
σ	0.54	0.58	0.58	[0.2 0.9]	[0.25 0.9]	[0.001, 0.0608]
α	0.45	0.50	0.42	[0.11 0.84]	[0.15 0.88]	[0.1 0.82]

z diagnostic table
True Model POMM, $N = 100$

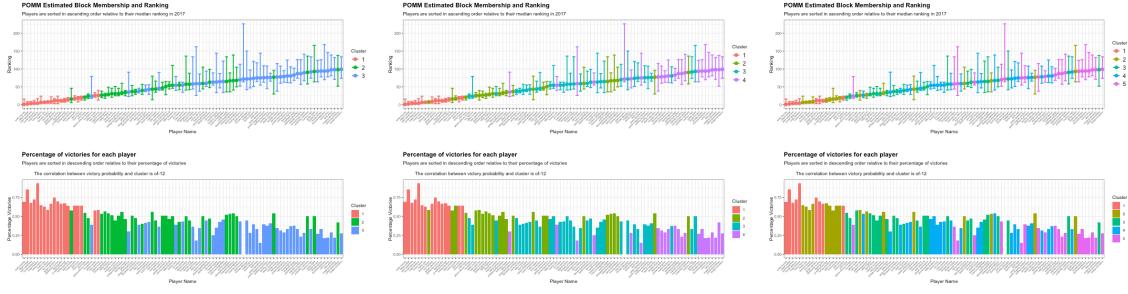
Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	11458.19	8062.46	11734.00	0.14	0.13	0.06	1.76	1.15	0.94	1.01	1.01	1.01
Simple	12846.46	8137.83	10373.67	0.10	0.18	0.05	1.63	1.17	0.92	1.01	1.07	1.01

P diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	947.67	1414.33	2051.6	0.11	0.08	0.05	29.96	29.59	32.31	1	1.00	1
Simple	1001.33	1365.83	1928.5	0.10	0.08	0.05	30.17	29.72	32.32	1	1.01	1

POMM hyperparameters diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
σ	3504	4202	5078	-0.01	0	0.01	37.25	34.56	34.16	1	1	1
α	10	13	15	0.96	0.96	0.97	25.42	25.07	24.88	1.19	1.01	1.1



13.2. Fixing $\sigma=0.01$.

$$\hat{P}^{POMM} = \begin{bmatrix} 0.50 & 0.65 & 0.79 \\ 0.34 & 0.50 & 0.64 \\ 0.21 & 0.35 & 0.50 \end{bmatrix}$$

$$\hat{P}^{POMM} = \begin{bmatrix} 0.50 & 0.58 & 0.68 & 0.76 \\ 0.42 & 0.50 & 0.57 & 0.67 \\ 0.32 & 0.43 & 0.50 & 0.57 \\ 0.24 & 0.33 & 0.43 & 0.50 \end{bmatrix}$$

$$\hat{P}^{POMM} = \begin{bmatrix} 0.50 & 0.58 & 0.67 & 0.72 & 0.79 \\ 0.42 & 0.50 & 0.57 & 0.67 & 0.71 \\ 0.33 & 0.43 & 0.50 & 0.57 & 0.67 \\ 0.28 & 0.33 & 0.43 & 0.50 & 0.58 \\ 0.22 & 0.29 & 0.33 & 0.42 & 0.50 \end{bmatrix}$$

z summary table
True Model POMM, $N = 100$

Method	WAIC		
	(a)	(b)	(c)
POMM model	-5220.111 21.17	-5181.395 19.32	-5233.632 20.07

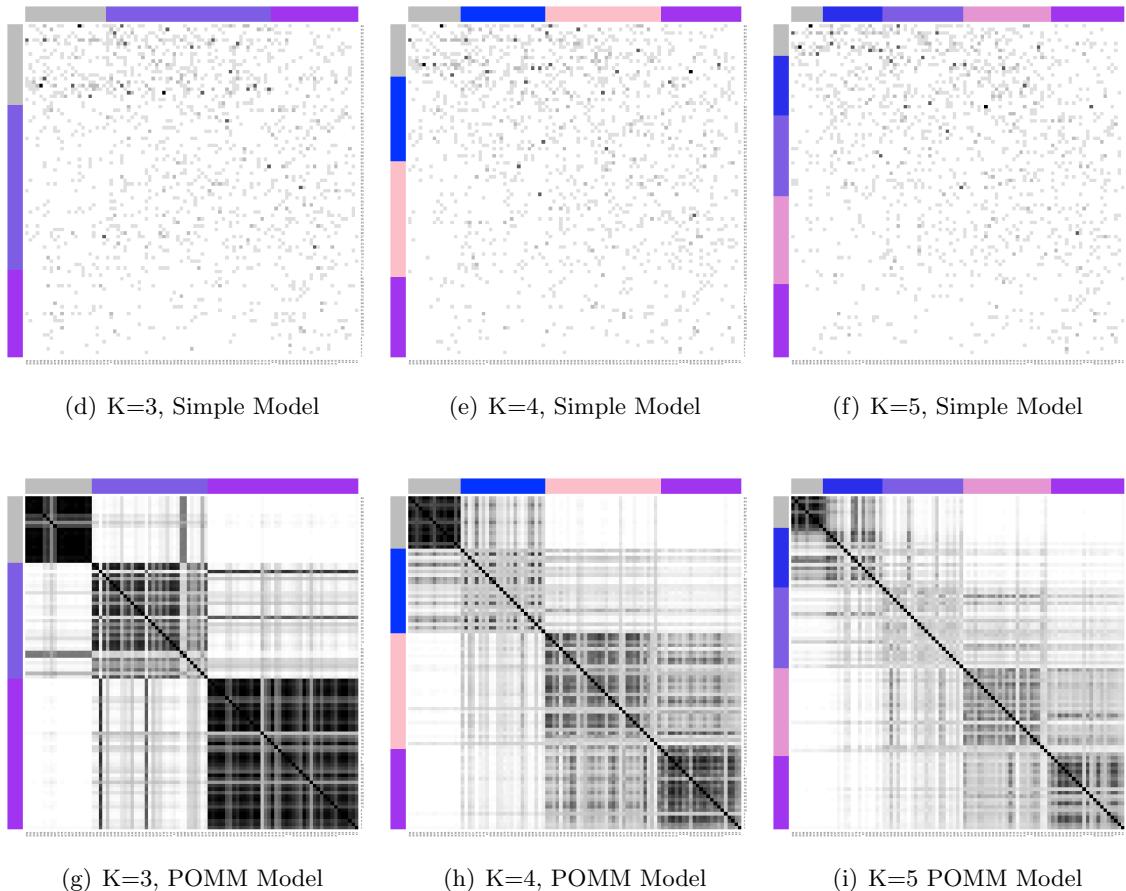


FIGURE 14. Adjacency Matrices simulated via the POMM Model

POMM Hyperparameters summary table
True Model POMM, $N = 100$

Method	$\hat{\theta}$			95% CI interval		
	(a)	(b)	(c)	(a)	(b)	(c)
σ	0.01	0.01	0.01	[0.01 0.01]	[0.01 0.01]	[0.01 0.01]
α	0.12	1.41	0.87	[0.1 0.15]	[0.1 2.74]	[0.1 1.8]

z diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	29353.04	20673.88	23681.54	0	0	0.01	8.93	15.05	16.53	1	1.13	1.16

P diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	2018.33	1900	2983.4	0.01	0.02	0.01	33.8	27.57	29.75	1	12.35	10.89

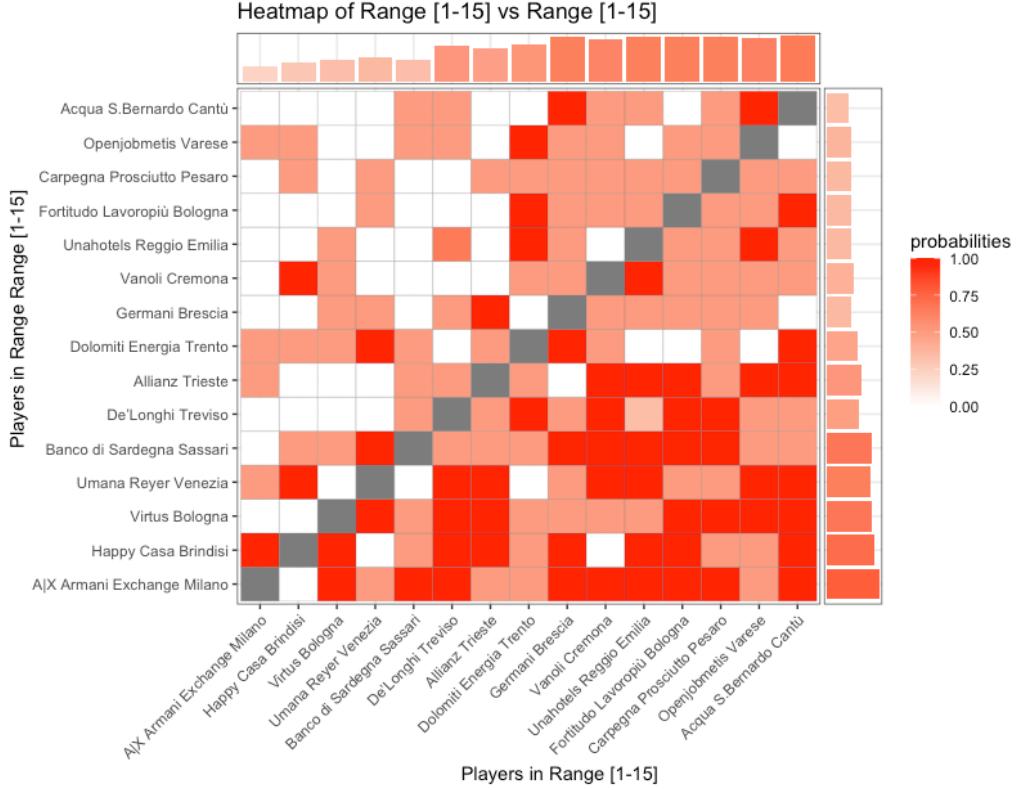


FIGURE 15. Heatmap filled with pairwise victory proportions. Players are arranged according to their final position within LBA ranking. On the side we can see their marginal victory proportions. We can see a clear increasing pattern in the marginals.

14. APPLICATION TO BASKET DATA

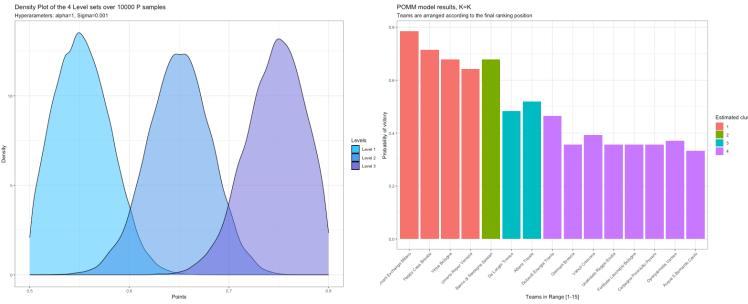
The principal league of Italian Basket is named Lega Serie A Basket, denoted from now on LBA.

We consider the season 2020/2021, which consists of 16 teams playing against each other 2 times. unfortunately, during that season, the Roma team filed for bankruptcy, and therefore we consider just 15 teams in our dataset.

We report here above a descriptive heatmap in figure (??).

The estimates below are obtained by fixing $\sigma^2 = 0.01$.

$$\hat{P}^{POMM} = \begin{bmatrix} 0.500 & 0.659 \\ 0.341 & 0.500 \end{bmatrix}$$



$$\hat{P}^{POMM} = \begin{bmatrix} 0.500 & 0.623 & 0.774 \\ 0.377 & 0.500 & 0.620 \\ 0.226 & 0.380 & 0.500 \end{bmatrix}$$

$$\hat{P}^{POMM} = \begin{bmatrix} 0.500 & 0.602 & 0.732 & 0.778 \\ 0.398 & 0.500 & 0.602 & 0.731 \\ 0.268 & 0.398 & 0.500 & 0.601 \\ 0.222 & 0.269 & 0.399 & 0.500 \end{bmatrix}$$

z summary table
True Model POMM, $N = 100$

Method	WAIC		
	(a)	(b)	(c)
POMM model	811.43 13.21	781.29 16.48	786.53 16.6

POMM Hyperparameters summary table
True Model POMM, $N = 100$

Method	$\hat{\theta}$			95% CI interval		
	(a)	(b)	(c)	(a)	(b)	(c)
σ	0.01	0.01	0.01	[0.01 0.01]	[0.01 0.01]	[0.01 0.01]
α	0.44	0.31	0.37	[0.1 0.84]	[0.1 0.65]	[0.1 0.71]

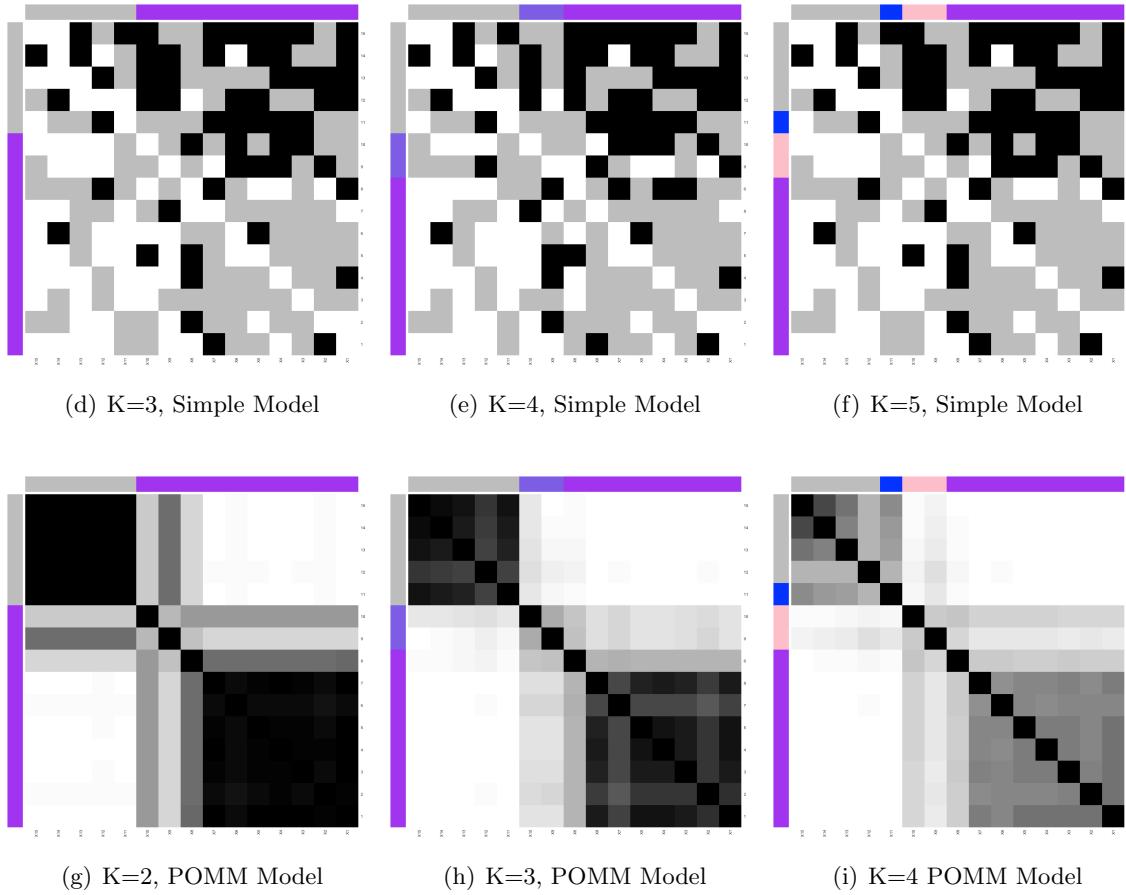


FIGURE 16. Adjacency Matrices simulated via the POMM Model

z diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	51953	34018	17349	0	0.02	0.11	42.56	5.58	17.9	1	1	1

P diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	2329	25	6.5	0	0.2	0.12	27.28	9.14	4.56	1	0.37	0.18

14.0.1. *Montecarlo algorithm.* First, it should be noted that within the same block we observe a substantial variability. We may have players that win against other players of the same cluster, against players of weaker clusters, or that simply do not win much. If we have players that substantially win against players of stronger clusters, it means that they are misclassified. So the fundamental variability driver is to be recognised in the blocks of the defeated players. Winning against Federer is not the same as winning against a newcomer. Those two victories should not be accounted in the same fashion.

However, one could argue that if the block exhibits a large amount of variability, this probably means that we should split the block further in two, to possibly account for different patterns of victories.

Another crucial point is data imbalance. Games are not drawn at random, which means that we will observe strongest player playing more with each other, since they are the ones remaining within the tournament for longer periods, and weaker players playing less against the strong ones and more among themselves.

The issue is that by assigning one cluster to a player is equivalent to equip them with a single probability to beat all the players within a given cluster.

15. APPLICATION TO MONKEY DATA

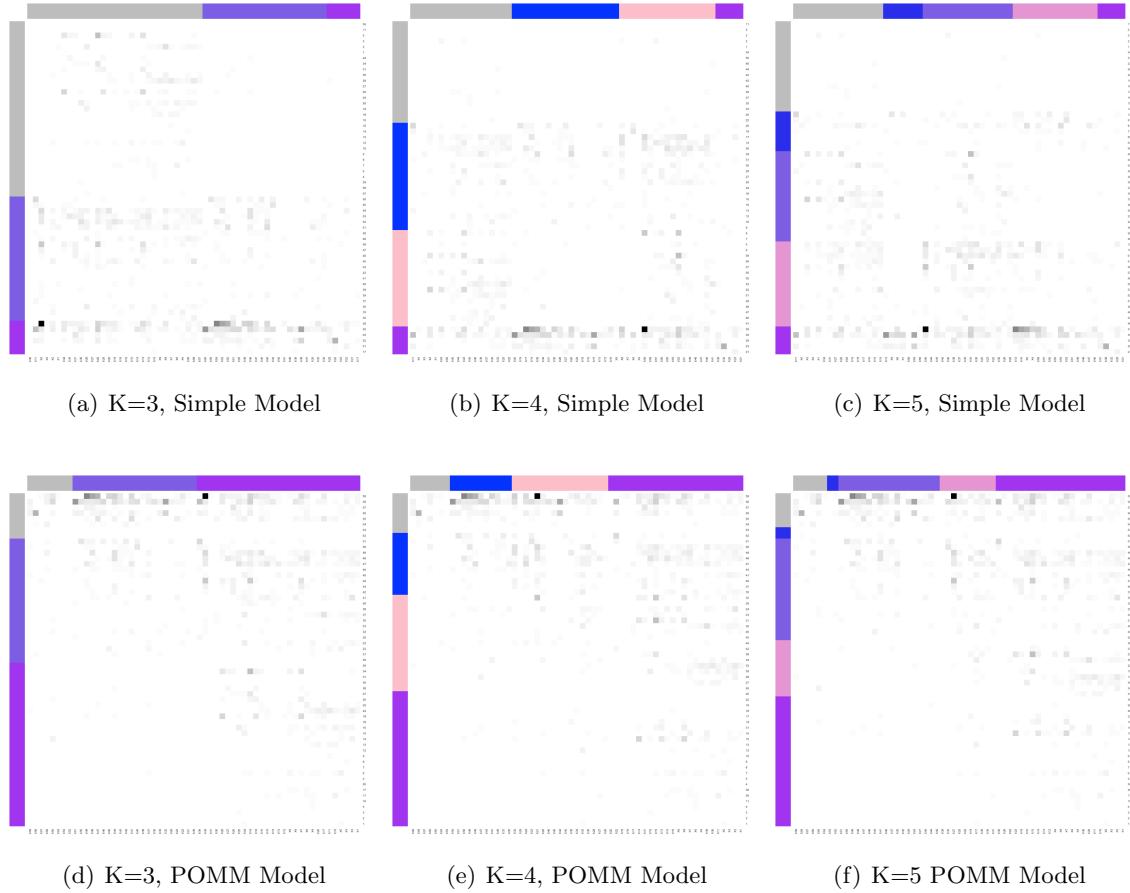


FIGURE 17. Adjacency Matrices simulated via the POMM Model

$$\hat{P}^{POMM} = \begin{bmatrix} 0.500 & 0.709 & 0.798 \\ 0.291 & 0.500 & 0.720 \\ 0.202 & 0.280 & 0.500 \end{bmatrix} \quad \hat{P}^{Simple} = \begin{bmatrix} 0.500 & 0.015 & 0.010 \\ 0.985 & 0.500 & 0.223 \\ 0.990 & 0.777 & 0.500 \end{bmatrix}$$

$$\hat{P}^{POMM} = \begin{bmatrix} 0.500 & 0.684 & 0.795 & 0.797 \\ 0.316 & 0.500 & 0.663 & 0.796 \\ 0.205 & 0.337 & 0.500 & 0.679 \\ 0.203 & 0.204 & 0.321 & 0.500 \end{bmatrix} \quad \hat{P}^{Simple} = \begin{bmatrix} 0.500 & 0.411 & 0.412 & 0.020 \\ 0.589 & 0.500 & 0.601 & 0.019 \\ 0.588 & 0.399 & 0.500 & 0.206 \\ 0.980 & 0.981 & 0.794 & 0.500 \end{bmatrix}$$

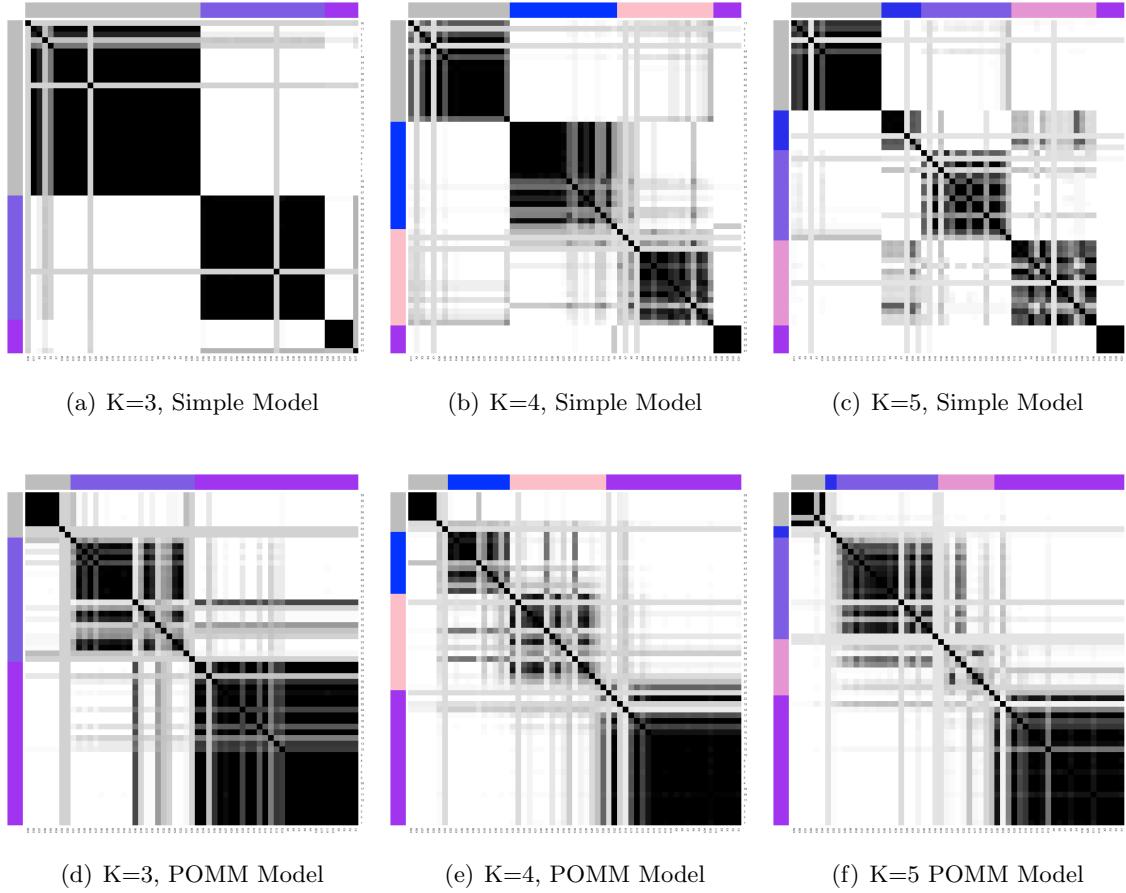


FIGURE 18. Adjacency Matrices simulated via the POMM Model

$$\hat{P}^{POMM} = \begin{bmatrix} 0.500 & 0.641 & 0.796 & 0.792 & 0.797 \\ 0.359 & 0.500 & 0.637 & 0.772 & 0.788 \\ 0.204 & 0.363 & 0.500 & 0.654 & 0.796 \\ 0.208 & 0.228 & 0.346 & 0.500 & 0.655 \\ 0.203 & 0.212 & 0.204 & 0.345 & 0.500 \end{bmatrix} \quad \hat{P}^{Simple} = \begin{bmatrix} 0.500 & 0.227 & 0.019 & 0.021 & 0.012 \\ 0.773 & 0.500 & 0.206 & 0.398 & 0.068 \\ 0.981 & 0.794 & 0.500 & 0.224 & 0.030 \\ 0.979 & 0.602 & 0.776 & 0.500 & 0.218 \\ 0.988 & 0.932 & 0.970 & 0.782 & 0.500 \end{bmatrix}$$

15.1. POMM model check.

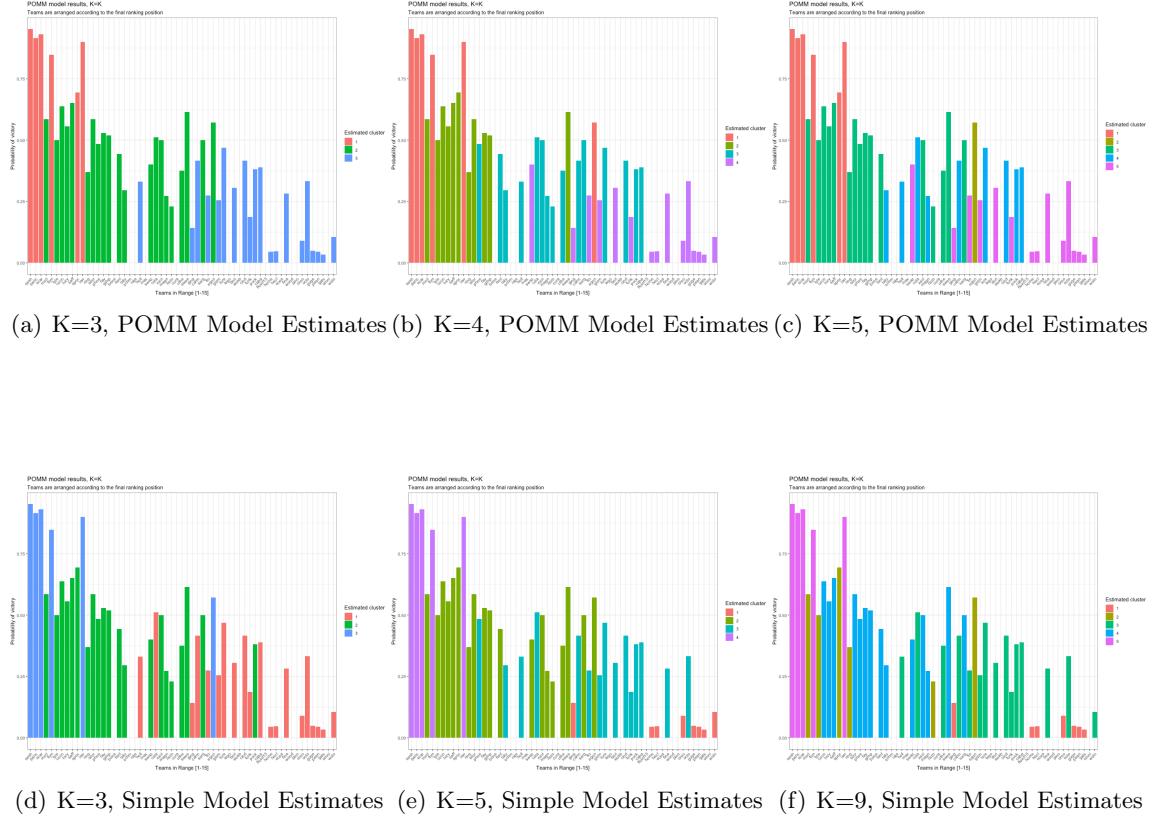


FIGURE 19. Co-Clustering Matrices obtained via the Simple Model (above) and the POMM model (below).

z summary table
True Model POMM, $N = 100$

Method	WAIC		
	(a)	(b)	(c)
POMM model	4705.28 148.72	4479.41 146.38	3990.14 107.89
Simple model	4475.26 167.55	3012.89 123.12	2973.91 121.29

POMM Hyperparameters summary table
True Model POMM, $N = 100$

Method	$\hat{\theta}$			95% CI interval		
	(a)	(b)	(c)	(a)	(b)	(c)
σ	0.54	0.58	0.58	[0.2 0.9]	[0.25 0.9]	[0.001, 0.0608]
α	0.45	0.50	0.42	[0.11 0.84]	[0.15 0.88]	[0.1 0.82]

z diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	36921.98	34213.27	27449.39	0.02	0.17	0.00	3.35	0.32	2.01	1	1	1.00
Simple	25819.39	20768.59	18504.95	0.00	0.00	0.03	37.60	0.01	0.13	Inf	Inf	12.48

P diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	1805	559.00	63.2	0.02	0.01	0.02	27.06	9.17	2.73	1.00	0.33	0.10
Simple	2177	905.67	240.9	0.02	0.00	0.00	27.76	10.95	3.06	128.46	0.53	4.12

POMM hyperparameters diagnostic table
True Model POMM, $N = 100$

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
σ	-	-	-	-	-	-	-	-	-	-	-	-
α	12	1466	1036	0.97	0.03	0.06	25.03	9.02	5.66	1.16	1	1

16. APPENDIX II: LIMITING CASE FOR σ

It is interesting to see what happens to the POMM prior when $\sigma \rightarrow \infty$. The idea is that, as the variance of the normals increase, the POMM prior should collapse on the Simple prior.

To prove this argument, first we check empirically and visually if the POMM prior distribution really converges to the Simple model distribution. Empirically we can run the Kolmogorov-Smirnov test to assess if there is a statistically significant difference between points $p_{ij} \sim POMM(\beta : \alpha; \sigma = \{0.01, 0.10, 0.50\}$ and the Simple prior $p_{ij} \sim Beta(1, 1)$. Then, we check visually the two distributions. Finally, we try to check analytically a convergence.

TABLE 21. Kolmogorov-Smirnov test
Data are generated via $p_{ij} \sim POMM(\alpha = 1, \beta_{max} = .8, \sigma)$

Method	p-value		
	$\sigma = 0.01$	$\sigma = 0.15$	$\sigma = 0.50$
POMM model	4.643e-09	3.998e-13	0.4163

In (??) we see that with values of σ equal to 0.5, we are unable to statistically distinguish points sampled from the Simple model and points sampled from the POMM model

$$(36) \quad \lim_{\sigma \rightarrow \infty} \prod_{k=1}^K \frac{1}{\sigma} \frac{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{p_{ij} - \mu(k)}{\sigma} \right)^2}}{\int_{-\infty}^{\beta} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{t - \mu(k)}{\sigma} \right)^2} dt - \int_{-\infty}^{0.5} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{t - \mu(k)}{\sigma} \right)^2} dt} =$$

$$(37) \quad \lim_{\sigma \rightarrow \infty} \prod_{k=1}^K \frac{1}{\sigma} \frac{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{p_{ij} - \frac{y(k) + y(k+1)}{2}}{\sigma} \right)^2}}{\int_{-\infty}^{\beta} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{t - \frac{y(k) + y(k+1)}{2}}{\sigma} \right)^2} dt - \int_{-\infty}^{0.5} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{t - \frac{y(k) + y(k+1)}{2}}{\sigma} \right)^2} dt}$$

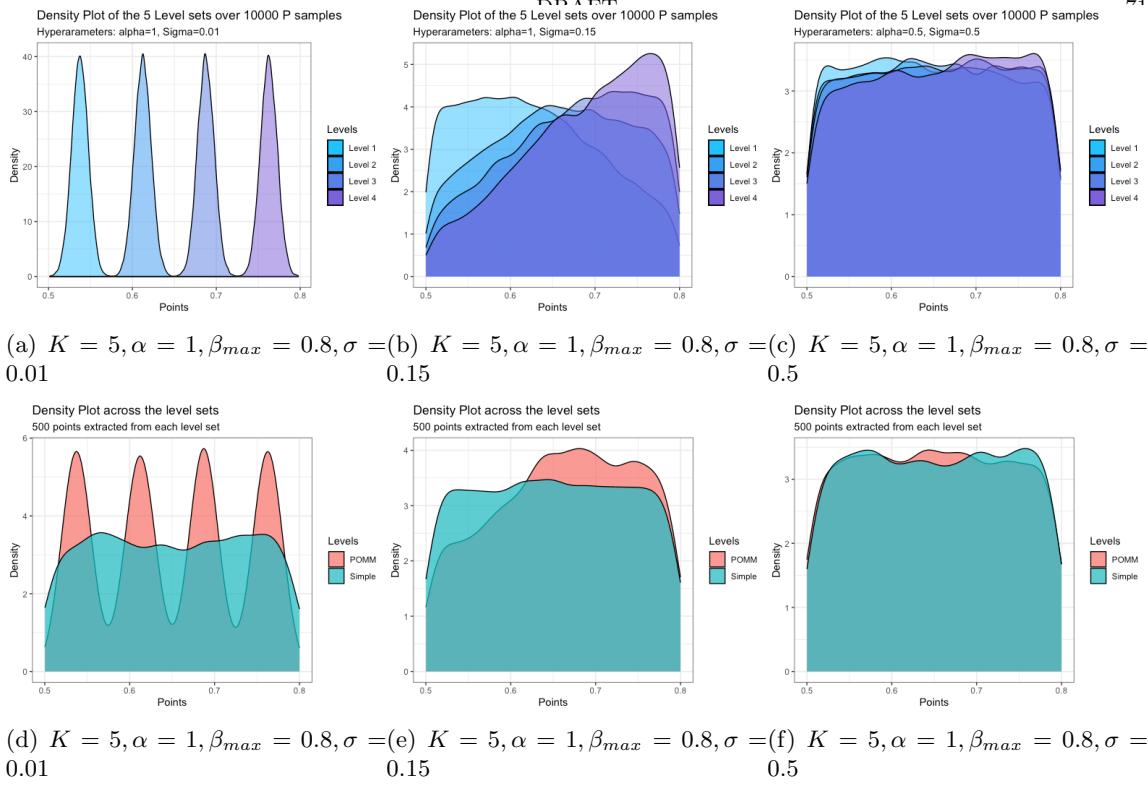


FIGURE 20. As σ increases, the distribution of the level sets of the POMM prior becomes indistinguishable from the distribution of the Simple model prior

If we substitute in the expression for $y_{(k)} = \left(\frac{(\beta_{\max} - 0.5)^{(1/\alpha)}}{K} \times k \right)^\alpha + 0.5$, we can simplify the following expression:

(38)

$$\begin{aligned} \frac{y_{(k)} + y_{(k+1)}}{2} &= \frac{1}{2} \left[\left(\frac{(\beta_{\max} - 0.5)^{(1/\alpha)}}{K} \times k \right)^\alpha + 0.5 + \left(\frac{(\beta_{\max} - 0.5)^{(1/\alpha)}}{K} \times (k+1) \right)^\alpha + 0.5 \right] \\ (39) \quad &= \frac{1}{2} \left[\left(\frac{(\beta_{\max} - 0.5)^{(1/\alpha)}}{K} \times k \right)^\alpha + \left(\frac{(\beta_{\max} - 0.5)^{(1/\alpha)}}{K} \times (k+1) \right)^\alpha + 1 \right] \end{aligned}$$

$$(40) \quad = \left[\frac{1}{2} \left(\frac{(\beta_{\max} - 0.5)}{K^\alpha} \times k^\alpha \right) + \left(\frac{(\beta_{\max} - 0.5)}{2K^\alpha} \times (k+1)^\alpha \right) + \frac{1}{2} \right]$$

$$(41) \quad = \frac{(\beta_{\max} - 0.5)}{2K^\alpha} \times (k^\alpha + (k+1)^\alpha) + \frac{1}{2}$$

Substituting back in the simplified expression:

(42)

$$\lim_{\sigma \rightarrow \infty} \prod_{k=1}^K \frac{1}{\sigma} - \frac{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{p_{ij} - \left[\frac{(\beta_{\max} - 0.5) \times (k^\alpha + (k+1)^\alpha) + \frac{1}{2}}{2K^\alpha} \right]}{\sigma} \right)^2}}}{\int_{-\infty}^{\beta} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{t - \left[\frac{(\beta_{\max} - 0.5) \times (k^\alpha + (k+1)^\alpha) + \frac{1}{2}}{2K^\alpha} \right]}{\sigma} \right)^2} dt - \int_{-\infty}^{0.5} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{t - \left[\frac{(\beta_{\max} - 0.5) \times (k^\alpha + (k+1)^\alpha) + \frac{1}{2}}{2K^\alpha} \right]}{\sigma} \right)^2} dt}$$

(43)

$$\lim_{\sigma \rightarrow \infty} \prod_{k=1}^K \text{TruncatedNormal}\left(\left[\frac{(\beta_{\max} - 0.5)}{2K^\alpha} \times (k^\alpha + (k+1)^\alpha) + \frac{1}{2}\right], \sigma^2; 0.5, \beta_{\max}\right)$$

Collapsing P parameter

We have that the likelihood, rewritten over blocks, takes the following expression:

$$f_Y(y | z, p) = \prod_{p=1}^k \prod_{q=1}^k \bar{\lambda}_{pq} p_{pq}^{\bar{y}_{pq}} (1 - p_{pq})^{\bar{m}_{pq}}$$

where

$$(44) \quad \bar{y}_{pq} = \sum_{i=1}^{N-1} \sum_{j=i+1}^N y_{ij} \mathbb{I}(z_i = p) \mathbb{I}(z_j = q)$$

$$(45) \quad \bar{m}_{pq} = \sum_{i=1}^{N-1} \sum_{j=i+1}^N (n_{ij} - y_{ij}) \mathbb{I}(z_i = p) \mathbb{I}(z_j = q)$$

$$(46) \quad \bar{\lambda}_{pq} = \sum_{i=1}^{N-1} \sum_{j=i+1}^N \binom{n_{ij}}{y_{ij}} \mathbb{I}(z_i = p) \mathbb{I}(z_j = q)$$

where \bar{y}_{pq} is the number of links between block p and block q , \bar{n}_{pq} is the number of missing links between the two blocks, and finally $\bar{\lambda}_{pq}$ computes the number of ways in which we could have obtained that number of links out of the total amount of possibilities.

Prior on P

The prior on the P entries takes the following expression

$$f_p(p_{pq} | \tilde{\alpha}_k, \tilde{\beta}_k, k) = \prod_{p=1}^K \prod_{q=1}^K \frac{1}{\text{Beta}(\tilde{\alpha}_k, \tilde{\beta}_k)} p_{pq}^{\tilde{\alpha}_k - 1} (1 - p_{pq})^{\tilde{\beta}_k - 1}$$

for $\tilde{\alpha}_k > 0, \tilde{\beta}_K > 0$, where for $q - p > 0$

$$(47) \quad \tilde{\alpha}_k = \left(\frac{1 - U_{(k)}}{\sigma^2} - \frac{1}{U_k} \right) \cdot \mu^2$$

$$(48) \quad \tilde{\mu}_k = \tilde{\alpha}_k \cdot \left(\frac{1}{U_{(k)}} - 1 \right)$$

so that

$$(49) \quad \mathbb{E} \left(p^{(l)} \right) = U_{(k)}$$

$$(50) \quad \text{Var} \left(p^{(l)} \right) = \sigma^2$$

where $p^{(l)} := \{p_{pq} : q - p = l \text{ for } l = 1, \dots, K - 1\}$, denoted as *Level Set*.

Instead, for $q - p = 0$, that is over $p^{(0)}$, we set $\tilde{\alpha}^K = \tilde{\beta}_k = 1$, so that the distribution of the level set $p^{(0)}$, that is, the diagonal of the P matrix is a uniform distribution with

$$(51) \quad \mathbb{E} \left[p^{(0)} \right] = 1/2$$

$$(52) \quad \text{Var} \left[p^{(0)} \right] = 1/12$$

The, using the beta-binomial conjugacy, we marginalize out P , obtaining

$$f_Y(Y | \mathbf{z}, \mathbf{U}_{(1), \dots, (K)}, \sigma^2, a) = \prod_{p=1}^K \prod_{q=1}^K \frac{B \left(\tilde{\alpha}_k + \bar{y}_{pq}, \tilde{\beta}_k + \bar{m}_{pq} \right)}{B \left(\tilde{\alpha}_k, \tilde{\beta}_k \right)}$$

where we use the notation $f_Y(y | z, \mathbf{U}_{(1), \dots, (K)}, \sigma^2, a)$ to explicit the dependency of $\tilde{\alpha}_k, \tilde{\beta}_k$ upon $\{\mathbf{U}_{(1), \dots, (K)}, \sigma^2, a\}$.

Collapsing the prior on γ

From this model we can write the joint distribution of $(\mathbf{z}, \boldsymbol{\gamma})$ as:

$$\begin{aligned} f(\mathbf{z}, \boldsymbol{\gamma}) &= f(\mathbf{z} | \boldsymbol{\gamma}) f(\boldsymbol{\gamma} | \boldsymbol{\alpha}) \\ &= \frac{1}{B(\boldsymbol{\alpha})} \gamma^{n_1} \gamma^{n_2} \dots \gamma^{n_K} \prod_{i=1}^K (\gamma_i^{\alpha_i - 1}) \\ (53) \quad &= \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^K (\gamma_i^{n_i + \alpha_i - 1}) \end{aligned}$$

where $n_k = \sum_{i=1}^N \mathbb{I}(z_i = k)$ is the number of items in block k . Now, we can marginalize out $\boldsymbol{\gamma}$ exploiting the Dirichlet-Categorial conjugacy

$$\begin{aligned}
f_{\mathbf{z}}(\mathbf{z} \mid \boldsymbol{\alpha}) &= \int p(\mathbf{z}, \boldsymbol{\gamma}) \boldsymbol{\gamma} \\
(54) \quad &= \frac{\sum_{i=1}^K \Gamma(\alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \int \gamma_i^{n_i + \alpha_i - 1} d\boldsymbol{\gamma} \\
&= \frac{\sum_{i=1}^K \Gamma(\alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \frac{\prod_{i=1}^K \Gamma(\alpha_i + n_i)}{\Gamma(\sum_{i=1}^K \alpha_i + n_i)} \\
&= \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\Gamma(\sum_{i=1}^K \alpha_i + n_i)} \prod_{i=1}^K \frac{\Gamma(\alpha_i + n_i)}{\Gamma(\alpha_i)} \\
(55) \quad &= \frac{\Gamma(A)}{\Gamma(A + N)} \prod_{i=1}^K \frac{\Gamma(\alpha_i + n_i)}{\Gamma(\alpha_i)}
\end{aligned}$$

where we have recognised in (??) the kernel of a Dirichlet($n_i + \alpha_i$) and where (??) is known as Dirichlet-Multinomial distribution and we denote it as $DM(\boldsymbol{\alpha})$

The full proportional posterior distribution

$$(56) \quad p(\mathbf{z}, \sigma^2, a, \bar{U}_{(s)\dots(k)}, K \mid YA) \propto$$

$$(57) \quad p(Y \mid \mathbf{z}, \sigma^2, a, \bar{U}_{(1)\dots(K)}, K) \cdot p(\mathbf{z} \mid A, K)$$

$$(58) \quad p(\sigma^2) \cdot p(\bar{U}_{(1)} \dots, (x) \mid a, K) \cdot p(a) \cdot p(K)$$

$$(59) \quad p(\mathbf{z} \mid YA) \propto p(Y \mid \mathbf{z}, \sigma^2, a) \cdot p(\mathbf{z} \mid A)$$

$$(60) \quad p(\sigma^2 \mid Y) \propto p(Y \mid \mathbf{z}, \sigma^2, a) \cdot p(\sigma^2)$$

$$(61) \quad p(a \mid Y) \propto p(Y \mid \mathbf{z}, \sigma^2) \cdot p(a)$$