

DRAFT

LAPO SANTI

CONTENTS

1. Introduction	2
2. Connection with Image Recognition	3
3. The SST model	5
3.1. Some considerations	6
4. The LST model	6
5. additional consideration	7
6. Modelling ordered blocks	7
7. Estimation	8
8. Point Estimate, Model Selection, and inference	10
9. Simulation Study from the Simple Model N=100	11
10. Simulation Study from the POMM Model N=100	15
10.1. POMM model check	15
11. Application to Tennis Data	18
11.1. POMM model check	19
11.2. Fixing $S=0.01$	22
12. Appendix I: Estimation Details	26
12.1. Updating z	26
12.2. Updating \mathbf{P}	26
13. Appendix II: POMM prior checks	27
13.1. Prior predictive check	27
13.2. MLE check	27

1. INTRODUCTION

When faced with a multitude of alternatives, individuals often strive to organize them into coherent blocks or groups to better understand the decision landscape. Furthermore, they aim to establish a meaningful order within these blocks, enabling them to prioritize alternatives based on preference. This process involves two fundamental tasks: block clustering and order-based ranking. While clustering involves categorizing alternatives into distinct blocks, ranking focuses on arranging these blocks in a specific order. These tasks are typically accomplished through human judgments, often in the form of pairwise comparisons.

In block clustering, the aim is to determine the inherent similarities among alternatives and group them accordingly. By comparing pairs of alternatives, individuals can identify common characteristics, shared attributes, or comparable features that contribute to their clustering. This process helps unveil the underlying structure of the alternatives, allowing decision-makers to comprehend the relationships and associations between them. Several techniques, such as hierarchical clustering and k-means clustering, have been employed to address this task effectively.

Conversely, in order-based ranking, the primary objective is to establish a preference-based order among the identified blocks of alternatives. By comparing pairs of blocks, individuals can discern the relative favorability of one block over another. These pairwise comparisons generate a ranking list that encapsulates the perceived preference or priority of each block. Various methodologies, including the Bradley-Terry model and pairwise comparison matrices, have been utilized to derive meaningful rankings from the collected preferences.

In this article, we propose a novel approach, termed Block Clustering and Order-based Ranking (BCOR), which unifies the tasks of block clustering and order-based ranking into a cohesive framework. The BCOR model introduces a dynamic parameter that governs the granularity of block clustering, allowing decision-makers to explore a spectrum of clustering options. By iteratively adjusting this parameter, the model can encompass a wide range of decision-making scenarios, from finely differentiated blocks to coarser groupings.

A key insight of the BCOR model lies in its ability to relate the number of blocks to the underlying ranking structure. As the number of blocks converges to the total number of alternatives, the model effectively transitions into a traditional ranking approach, providing a complete ordering of the alternatives. Conversely, by intentionally reducing the number of blocks, decision-makers are presented with distinct groups of choices, each requiring preference considerations within its own subset. This approach offers a nuanced perspective on decision-making, allowing individuals to differentiate between highly favored groups and those that are comparatively less preferred.

The BCOR model provides a flexible and adaptive solution for organizing and prioritizing alternatives in various decision-making contexts. Its application extends beyond conventional clustering and ranking tasks, empowering decision-makers to explore the continuum between comprehensive rankings and granular groupings. Additionally, the model can be

tailored to incorporate different types of pairwise comparisons, enabling its utilization in diverse domains and decision scenarios.

To evaluate the effectiveness of the BCOR model, we conducted experiments using real-world datasets encompassing a wide range of decision contexts. The results demonstrate the model’s ability to generate meaningful block clusters and order-based rankings, outperforming traditional approaches that solely focus on clustering or ranking tasks.

The contributions of this work can be summarized as follows:

We introduce the novel problem of block clustering and order-based ranking, bridging the gap between these two fundamental decision tasks. We propose the BCOR model, which provides a unified framework to accommodate various levels of granularity in decision-making, from complete rankings to distinct preference-based groups. We showcase the versatility of the BCOR model through experiments on real-world datasets, highlighting its superior performance compared to existing methods. By integrating block clustering and order-based ranking, the BCOR model offers decision-makers a comprehensive tool to navigate complex decision landscapes

2. THE SIMPLE MODEL

2.1. Simple 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

$$\begin{aligned}
 (1) \quad p(y|z, P, K) &= \prod_{i=2}^{N-1} \prod_{j=i}^N p(y_{ij}|z, P, K) \\
 (2) \quad &= \prod_{i=2}^{N-1} \prod_{j=i}^N \binom{n_{ij}}{y_{ij}} p_{z_i, z_j}^{y_{ij}} (1 - p_{z_i, z_j})^{n_{ij} - y_{ij}}
 \end{aligned}$$

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 .

2.2. 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 $\boldsymbol{\theta}$, following the common practice in the literature, we can express the marginal distribution of z as:

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

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 $K > 0$.

3. 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:

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

Instead, under SST, we would observe:

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

With regard of *LST*, it is a generalisation of the two axioms and therefore it includes both aforementioned cases (3) and (4), 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.

4. THE SST MODEL

We consider a Partially Ordered Markov Model (POMM) for a set of entries $p_{ij} \in L^{(k)}$, where $L^{(k)}$ represents a level set. These entries are assumed to be identically and independently distributed according to a truncated normal distribution:

$$p_{ij} \mid (y^{(k)}, y^{(k+1)}) \sim \text{TruncatedNormal}(\mu^{(k)}, \sigma^{2(k)}; 0.5, \beta_{\max})$$

Here:

- $\mu^{(k)}$ is the mean, which corresponds to the midpoint of the level set $L^{(k)}$, defined as $\mu^{(k)} = \frac{y^{(k)} + y^{(k+1)}}{2}$.
- σ^2 is the variance parameter constant across the level sets $L^{(1)}, \dots, L^{(K)}$.
- 0.5 and β_{\max} are the lower and upper truncation bounds.

We also place a uniform hyper-prior on the parameter α :

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

which in turn, together with β_{\max} , it provides the truncations of the level sets as follows:

$$(11) \quad y^{(k)} = \left(\frac{(\beta_{\max} - 0.5)^{(1/\alpha)}}{K} \times k \right)^\alpha + 0.5 \quad \text{for } k = 0, \dots, K$$

We also place a uniform hyper-prior on the parameter σ^2

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

which is the variance of the level sets truncated normal distribution.

Supposition: When $\sigma^2 \rightarrow \infty$, the joint distribution of the level sets collapses to a uniform and therefore we are back to a uniform distribution.

4.1. Some considerations. This model is very flexible, since for a given number of clusters K , it can accommodate for different means across the level sets by changing the slope of α . For an $\alpha \in [0, 1]$, we will have the clusters concentrated towards β_{\max} , meaning that the probabilities of victory in the tournament are clear-cut, symptom of an higher predictability overall. Instead, if $\alpha \in [1, 3]$, this means that the probabilities are more concentrated towards $1/2$, meaning the chances of victory are not so much clear.

However, the issue is that the number of clusters is not fixed. We would like the model to select by itself the number of clusters and ideally also to place the mean of the distribution wherever it prefers.

One question may be: what is the difference between a model with a very high α , meaning with the probabilities concentrated towards $1/2$ but with three clusters, and a model with just one cluster, and all the players having the same mean?

Maybe the parameter α is redundant and we should find a clever way so that, when we add or remove a cluster, we assign a mean of victory to that given level set.

However, the means of the level sets μ_1, \dots, μ_K should be ordered as well and we should put a constraint on them.

Another point that I have is the following:

5. THE LST MODEL

The latent skill variables for clusters are denoted by s_c for cluster c . We model these latent skills using a Gaussian Process (GP) prior:

$$s_c \sim \mathcal{GP}(0, k(\cdot, \cdot))$$

Here, \mathcal{GP} is the Gaussian Process distribution, and $k(\cdot, \cdot)$ is the covariance (kernel) function that captures the relationships between clusters.

The victory probability for a player from cluster A against a player from cluster B is modeled using a logistic function f that takes the difference in latent skills into account:

$$p_{A,B} = f(s_A - s_B) = \frac{1}{1 + \exp(-\beta \cdot (s_A - s_B))}$$

Where β is a parameter that controls the slope of the logistic function and affects how fast the probability changes with the difference in skills.

The likelihood of observing y_{ij} victories out of n_{ij} matches between clusters i and j is given by a binomial distribution:

$$y_{ij} \sim \text{Binomial}(n_{ij}, p_{i,j})$$

The complete model, combining the Gaussian Process prior on latent skills and the logistic function for victory probabilities, can be written as:

For each cluster c :

$$s_c \sim \mathcal{GP}(0, k(\cdot, \cdot))$$

For each match between clusters i and j :

$$p_{i,j} = \frac{1}{1 + \exp(-\beta \cdot (s_i - s_j))}$$

$$y_{ij} \sim \text{Binomial}(n_{ij}, p_{i,j})$$

6. ADDITIONAL CONSIDERATION

The LST and the SST and the WST models can all of them being expressed as part of the same class of models.

7. ESTIMATION

For the moment, we want to infer just $\theta = \{z, P, \alpha, S\}$, meaning that we treat K as a known constant. The estimation strategy is a Hybrid MCMC algorithm. Since simulating from the conditional distribution $p(\theta_i | \theta_j, j \neq i)$ is unfeasible or computationally expensive, we substitute the simulation from the full conditional distribution with a simulation from a proposal distribution q_i . Referencing Muller's (1991) work, the Hybrid modification is as follows:

Algorithm 1 Metropolis-within-Gibbs MCMC

```

for  $i = 1, \dots, p$  given  $(\theta_1^{(t+1)}, \dots, \theta_{i-1}^{(t+1)}, \theta_i^{(t)}, \dots, \theta_p^{(t)})$  do
  1. Simulate
  (12)  $\theta'_i \sim q_i(\theta_i^{(t)} | \theta_1^{(t+1)}, \dots, \theta_i^{(t)}, \theta_{i+1}^{(t)}, \dots, \theta_p^{(t)})$ 
  2. Take
  (13)  $\theta_i^{(t+1)} = \begin{cases} \theta_i^{(t)} & \text{with probability } 1 - r_i, \\ \theta'_i & \text{with probability } r_i, \end{cases}$ 
  where
  (14)  $r_i = 1 \wedge \left\{ \frac{p(\theta'_i | \theta_i^{(t)} | \theta_1^{(t+1)}, \dots, \theta_i^{(t)}, \theta_{i+1}^{(t)}, \dots, \theta_p^{(t)})}{p(\theta_i^{(t)} | \theta_i^{(t)} | \theta_1^{(t+1)}, \dots, \theta_i^{(t)}, \theta_{i+1}^{(t)}, \dots, \theta_p^{(t)})} \right\}$ 
end for

```

7.0.1. *Adaptive algorithm for $\theta = \{P, \alpha, S\}$.* We specify the proposal distributions in (7) above as

$$\theta'_i \sim \text{Normal}(\theta_i^{(t-1)}, \sigma_{\theta_i}^2)$$

whose sampled value is accepted or rejected by evaluating the logarithm of (9). Choosing a correct $\sigma_{\theta_i}^2$ 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 $K(K-1)/2+2$ parameters i ($1 \leq i \leq K(K-1)/2+2$), 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 Sherlock 2009. 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.

→ Insert here plots of convergence to the acceptance ratio

We specify in the Appendix the full expression for the ratio of $\theta = \{P, \alpha, S\}$ in (9).

7.0.2. *Adaptive Algorithm for $\theta = z$.* When dealing with $\theta = z$, a discrete parameter, we need to adapt the formulation while maintaining the underlying concept. In the case of the POMM model, the labels $k = 1, \dots, K$ are ordered, and therefore, we can define a distance metric between these labels. Let us denote the distance between k and k' as $d(k, k')$, which can be expressed as:

$$(15) \quad d(k, k') = |k - k'|$$

If the acceptance rate for a particular player 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. To achieve this, we assign a sampling probability to each label that is inversely related to its distance from the current label. Specifically, we define $p(k') = p(|k' - k|) = \text{Normal}(0, \sigma_i^2)$, where σ_i^2 is adapted as above. A larger variance assigns higher probabilities to distant labels, while a smaller variance favors closer labels. Finally, we employ a multinomial distribution to sample the next label k' :

$$(16) \quad k' \sim \text{Multinomial}(1, K, p(|k' - k|))$$

By using this approach, we can adapt the algorithm to explore labels based on their distances from the current label.

We specify in the Appendix the full expression for the ratio of $\theta = \{z\}$ in (9).

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 three main measures:

- Variation of Information (VI): to fully utilise this posterior and engage in inference directly within the partition space, we adopt the decision-theoretic approach introduced by Wade and Ghahramani (2018) for block modeling. 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 $\log 2V$, where V 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. Refer to Wade and Ghahramani (2018) for a detailed exploration of the key properties of the vi metric. Within this framework, a formal Bayesian point estimate for z is obtained by selecting the partition with the lowest averaged vi distance from the other clusterings
- WAIC: While the WAIC yields practical and theoretical advantages and has direct connections with Bayesian leave-one-out cross-validation, thus providing a measure of edge predictive accuracy, the calculation of the WAIC only requires posterior samples of the log-likelihoods for the edges: $\log p(y_{ij}|z, P, \alpha) = y_{ij} \log p_{z_i, z_j} + (n_{ij} - y_{ij}) \log(1 - p_{z_i, z_j})$, $i = 2, \dots, N, j = 1, \dots, i - 1$.
- Misclassification error: predicting the group membership z_{N+1} of a new player may also be of interest. We can derive the estimate of the block probabilities for new players based on their early matches with some of the existing players.

$$(17) \quad \begin{aligned} p(z_{N+1} = k | \mathbf{Y}, y_{N+1}, \hat{z}) &\propto p(y_{N+1} | \mathbf{Y}, \hat{z}, z_{N+1} = k) \cdot p(z_{N+1} = k | \hat{z}) \\ &= p(y_{N+1} | \hat{z}, z_{N+1} = k) \cdot p(z_{N+1} = k | \hat{z}) \end{aligned}$$

where $p(z_{N+1} = k | \mathbf{Y}, y_{N+1}, \hat{z})$ is the posterior probability of the new node $N + 1$ to belong to the block k , given the previously observed data Y , the new node's data y_{N+1} and the estimated labels \hat{z} . On the right hand side of the expression above, $p(y_{N+1} | \mathbf{Y}, \hat{z}, z_{N+1} = k)$ represents the likelihood of observing y_{N+1} given the previously observed data Y and the estimated labels \hat{z} , which, due to conditional independence, is the same as conditioning just on \hat{z} . Finally, $p(z_{N+1} = k | \hat{z})$ represents the prior probability of label k for the new node $N + 1$ given \hat{z} , which we can approximate with the relative size of the blocks n_k .

9. SIMULATION STUDY FROM THE SIMPLE 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 Simple Model. We want to compare how it performs compared to the POMM extension and other state-of-the-art alternatives. Each tournament had a different number of blocks in the underlying structure. We set the total number of games $M := 0.5 * \sum_{i,j}^N n_{ij} = \sum_{i,j}^N y_{ij} = 4000$, which is the average number of matches played in one year of tennis tournaments. We divided the players into three, five and nine blocks ($K = 3, 5, 9$ respectively). In Figure (1), we display the three simulated tournaments, where the difficulty of accurately determining the group membership increases as the number of games increases with the number of blocks.

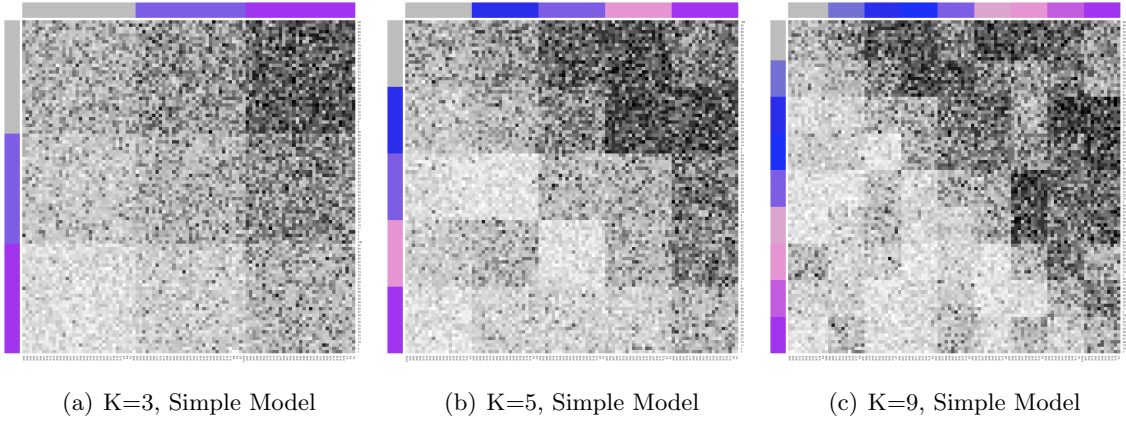


FIGURE 1. Adjacency Matrices simulated via the Simple Model

We compare the performance of the Simple model with the POMM one. We fixed arbitrarily $\beta_{\max} = .75$. In table (??) we report the results of the simulation. In the three cases, for the Simple and the POMM model, we compare the WAIC, the VI distance and the misclassification error, obtained by considering 100 new incoming players which get to play just with 10 players each. We also compare the labels estimated against the regularised spectral clustering algorithm and the Louvain algorithm. The Simple model is the best performing relative to the other three alternatives.

In figure (2) we report the estimated co-clustering matrices resulting from the simulation process.

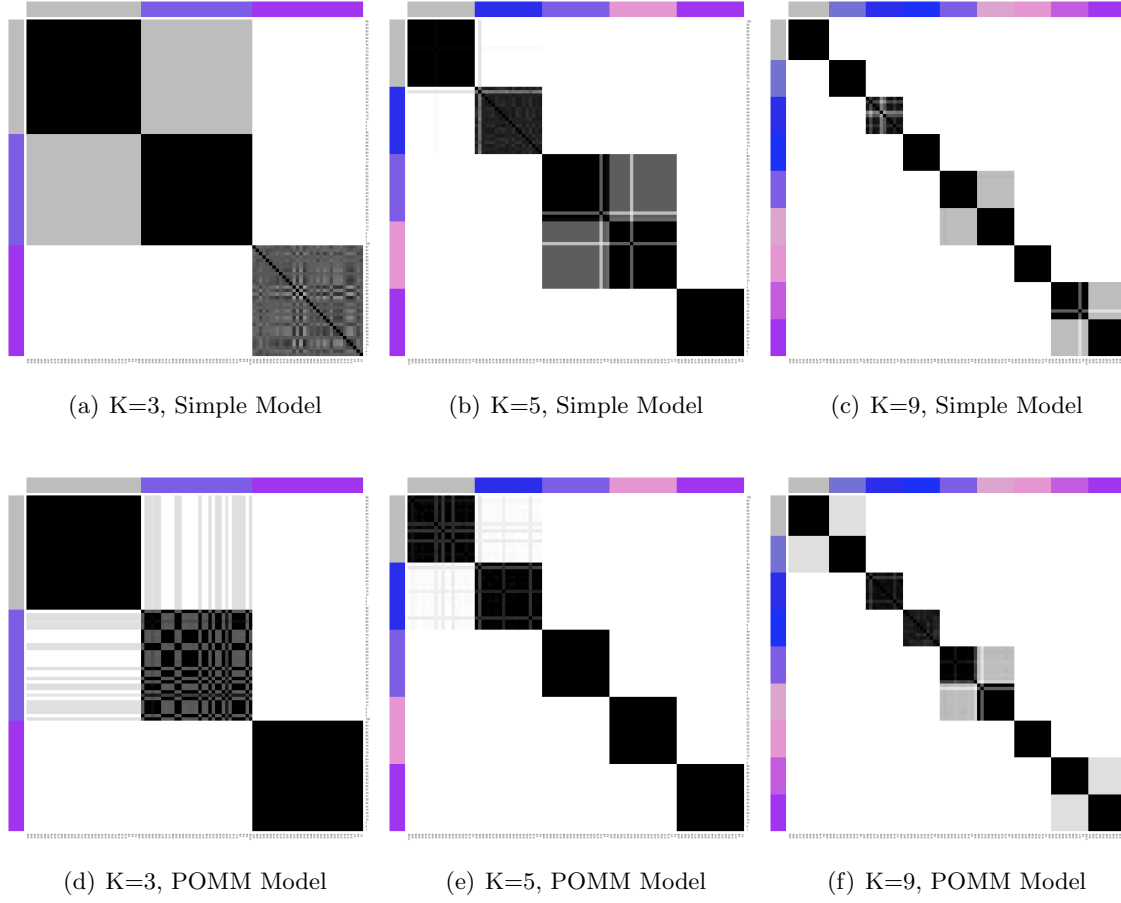


FIGURE 2. Co-Clustering Matrices obtained via the Simple Model(above) and the POMM model (below).

P summary table
True Model Simple, $N = 100$

Fitted Model	\overline{MAE}			% within-95% CI interval			\overline{CI} interval length		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0.03	0.06	0.07	100	40	63.89	0.21	0.04	0.15
Simple model	0.04	0.08	0.06	100	60	66.67	0.10	0.14	0.16

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

Method	VI distance _{MAP}			VI distance _{VI lb}			WAIC		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0	0	0.00	0.00	0.0	0.22	-6380.08 16.90	-6242.73 19.03	6358.67 18.65
Simple model	0	0	0.27	0.67	0.40	0.44	-6334.84 17.60	-6184.29 19.04	-6356.23 18.61

POMM hyperparameters summary table
True Model Simple, $N = 100$

Fitted Model	$\hat{\theta}$			95% CI interval		
	(a)	(b)	(c)	(a)	(b)	(c)
S	0.29	0.15	0.27	[0.01 0.83]	[0.04 0.55]	[0.07 0.78]
α	0.55	0.36	0.37	[0.18 0.89]	[0.1 0.61]	[0.1 0.75]

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

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	374.75	719.47	862.00	0.09	0.01	0.26	0.01	0.01	31.26	Inf	6.26	1.42
Simple	1926.61	533.16	20.64	0.03	0.21	0.06	0.01	0.01	0.01	Inf	Inf	Inf

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

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	1520	1761.6	1419.61	0.09	0.00	0.07	25.12	35.37	32.42	46.75	4.27	7.96
Simple	1265	1545.3	1493.03	0.07	0.09	0.07	30.80	33.01	32.58	14.98	12.00	8.05

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

Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
S	1362	435	0	0.32	0.4	0	33.35	29.64	0	2.5	1.58	1.54
α	11	16	19	0.93	0.91	0.9	24.66	25.12	24.61	1.11	1.04	1.93

9.0.1. *Simple Model check.*

10. 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 Simple model, now we are doing the same, with similar parameters ($K = 3, 5, 9, M = 4000$ and $\beta_{\max} = .75$). Here are the results.

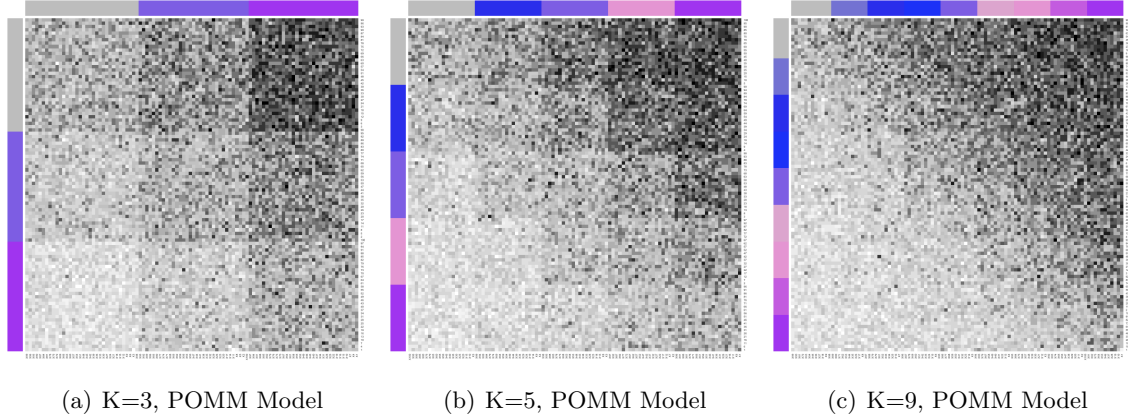


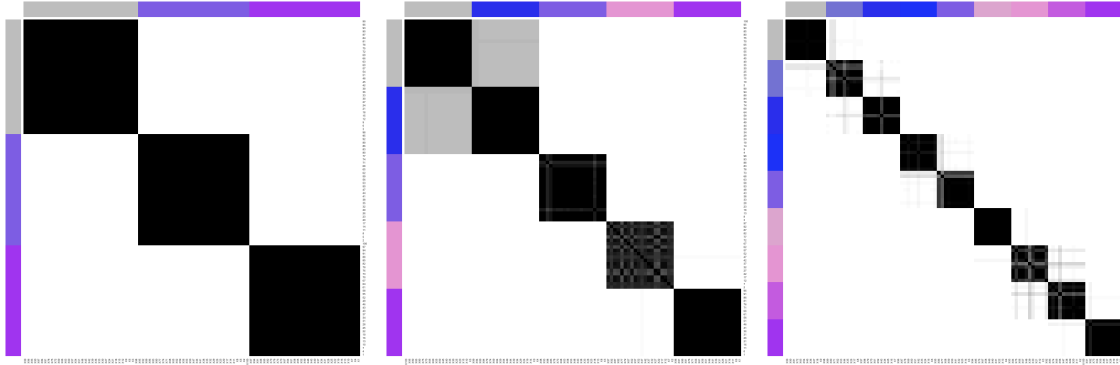
FIGURE 3. Adjacency Matrices simulated via the POMM Model

In table (??) we report the results of the simulation. As before, for the Simple and the POMM model, we compare the WAIC, the VI distance and the misclassification error, obtained $N_{new} = 100$. Also here we compare clustering performance against that of the regularised spectral clustering algorithm and the Louvain algorithm. The POMM model is the best performing relative to the other three alternatives.

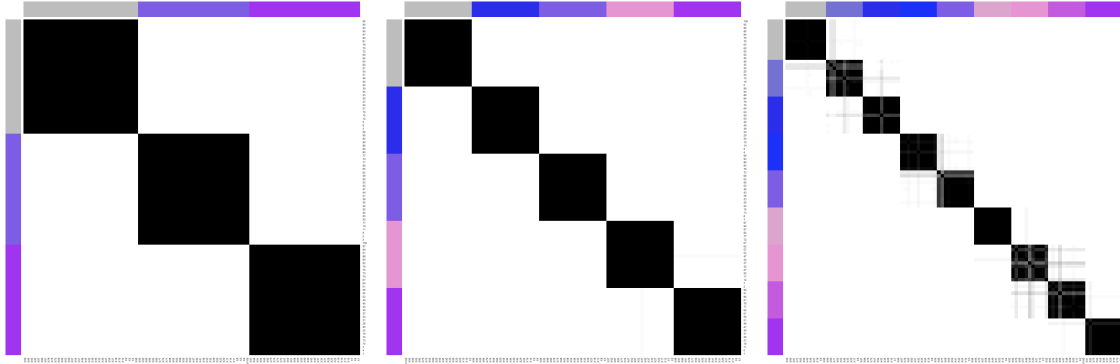
P summary table
True Model POMM, $K = 3$, $N = 100$

Fitted Model	\overline{MAE}			% within-95%-CI interval			$\overline{\text{CI interval length}}$		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0.03	0.02	0.01	100%	100%	94.44%	0.13	0.08	0.06
Simple model	0.00	0.05	0.01	100%	90%	97.22%	0.01	0.13	0.07

10.1. POMM model check.



(a) K=3, Simple Model Estimates (b) K=5, Simple Model Estimates (c) K=9, Simple Model Estimates



(d) K=3, POMM Model Estimates (e) K=5, POMM Model Estimates (f) K=9, POMM Model Estimates

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

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

Method	VI distance _{MAP}			VI distance _{VI lb}			WAIC		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM model	0	0	0	0	0.0	0	-6496.08 17.55	-6620.84 17.97	6658.97 17.78
Simple model	0	0	0	0	0.4	0.0	-6553.32 17.60	-6566.43 17.84	6662.66 17.96

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)	(a)	(b)	(c)
S	0.1487	0.084	0.0196	[0.0033 0.7575]	[4e-04 0.5602]	[0.0012 0.0625]	0.01	0.01	0.01	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	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	0	15.19	2167.05	0.06	0.20	0.16	0.01	0.01	0.05	Inf	Inf	8.53
Simple	0	670.73	2901.09	-0.01	0.04	0.02	0.01	0.01	0.03	NA	Inf	6.93

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	1821.67	1296.1	1035.03	0.05	0.06	0.06	27.13	33.31	32.94	32.56	10.94	4.45
Simple	2121.33	1505.6	1350.92	0.00	0.04	0.01	29.88	33.29	32.70	1.00	12.26	3.26

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)
S	938	548	316	0.1	0.27	0.54	25.42	12.03	6.71	3.5	2.15	8.15
α	33	341	180	0.82	0.3	0.41	23.83	17.82	13.81	1.07	1.44	2.41

11. APPLICATION TO TENNIS DATA

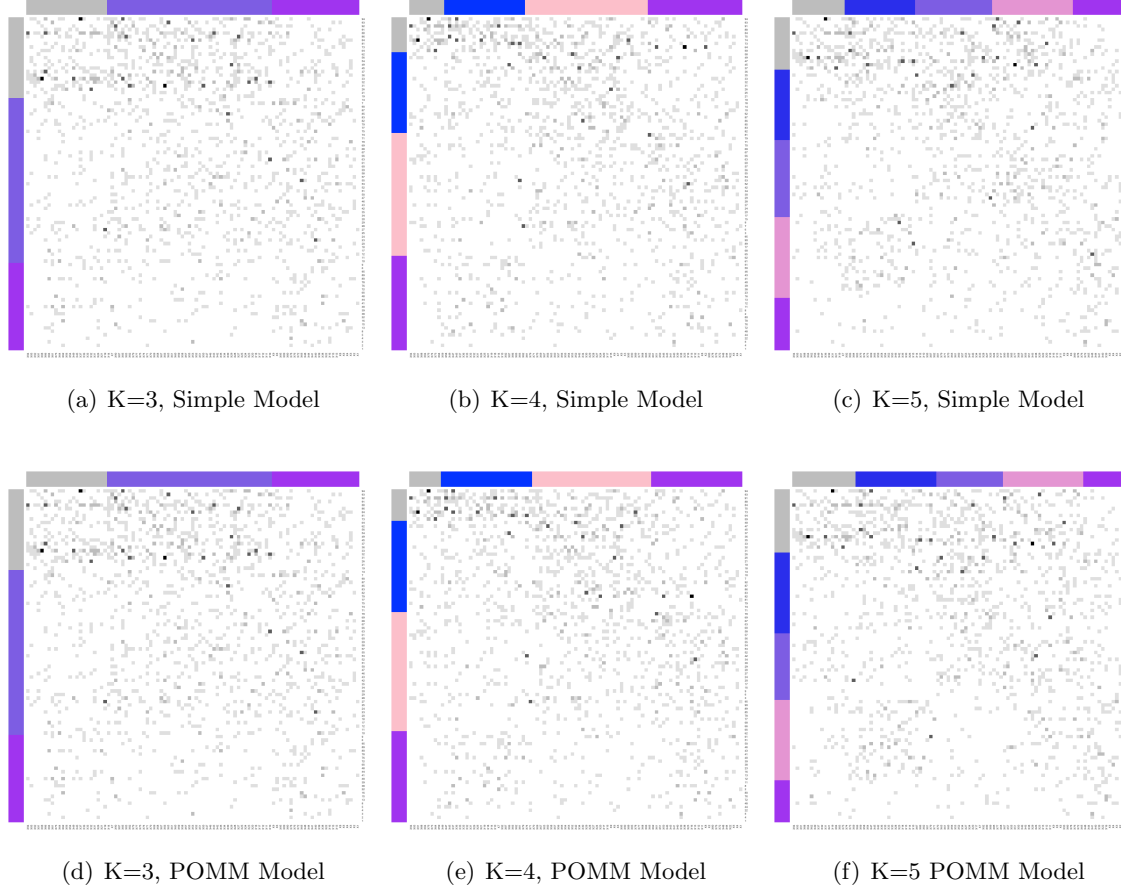


FIGURE 5. Adjacency Matrices simulated via the POMM Model

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

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

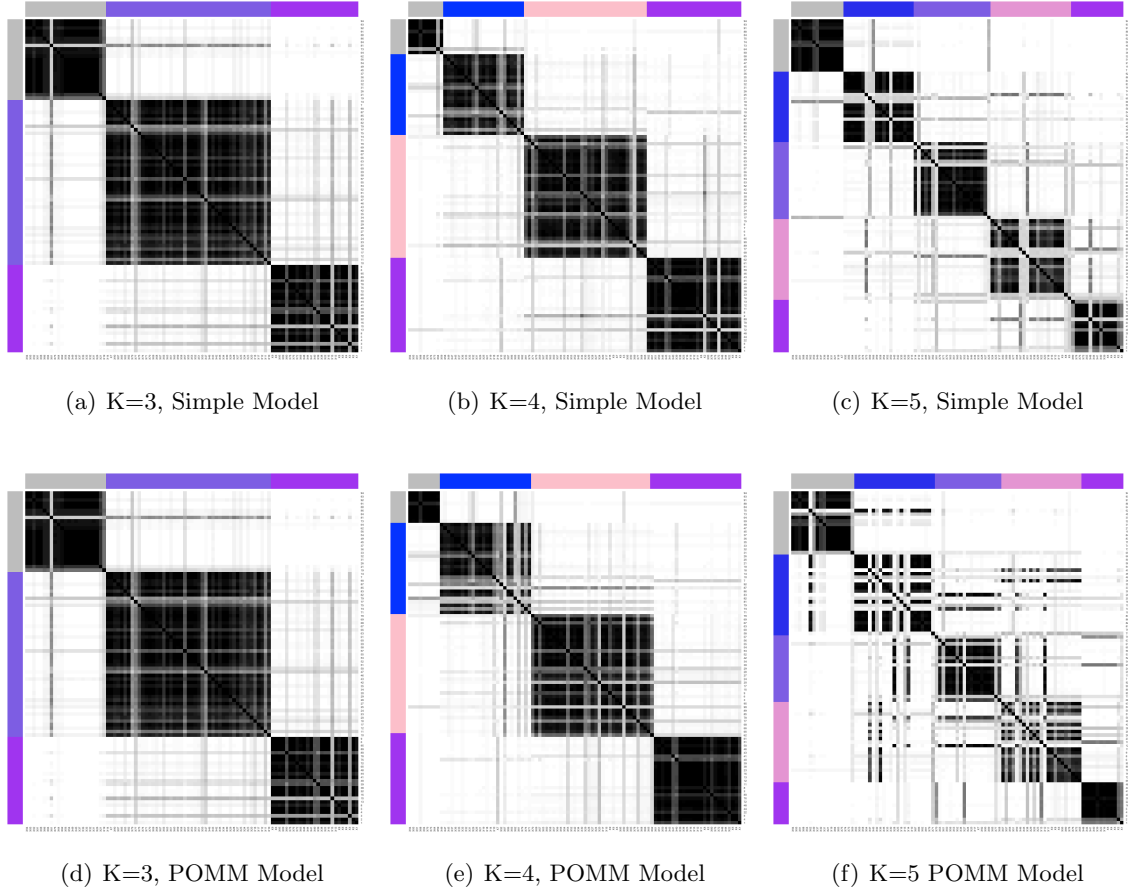
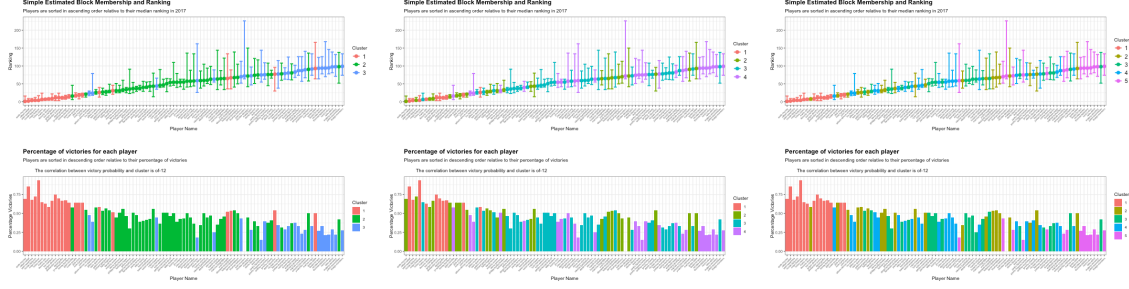


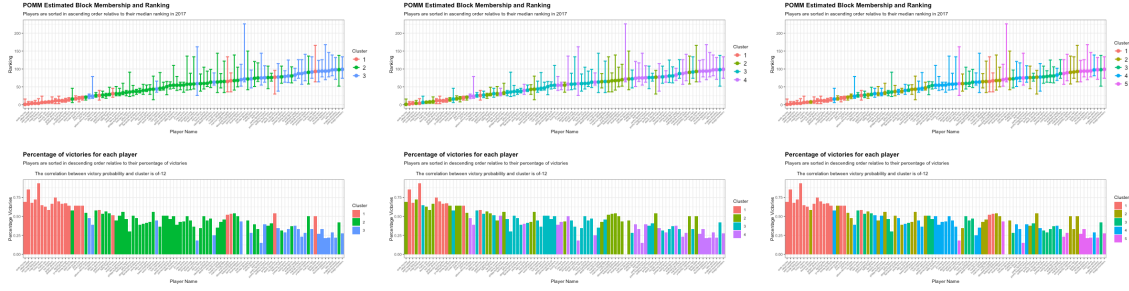
FIGURE 6. Adjacency Matrices simulated via the POMM Model

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

11.1. POMM model check.



(a) K=3, Simple Model Estimates (b) K=5, Simple Model Estimates (c) K=9, Simple Model Estimates



(d) K=3, POMM Model Estimates (e) K=5, POMM Model Estimates (f) K=9, POMM Model Estimates

FIGURE 7. 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)
S	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, $N = 100$

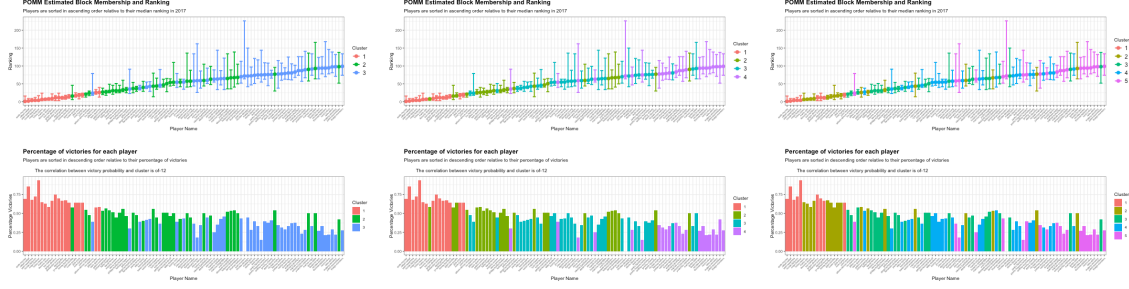
Fitted Model	ESS			ACF ₃₀			% accepted			Gelman-Rubin		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
POMM	11458.19	8062.46	11734.00	0.14	0.13	0.06	1.76	1.15	0.94	1.01	1.01	1.01
Simple	12846.46	8137.83	10373.67	0.10	0.18	0.05	1.63	1.17	0.92	1.01	1.07	1.01

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)
S	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



(a) K=3, POMM Model Estimates (b) K=5, POMM Model Estimates (c) K=9, POMM Model Estimates

11.2. Fixing $S=0.01$.

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

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

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

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

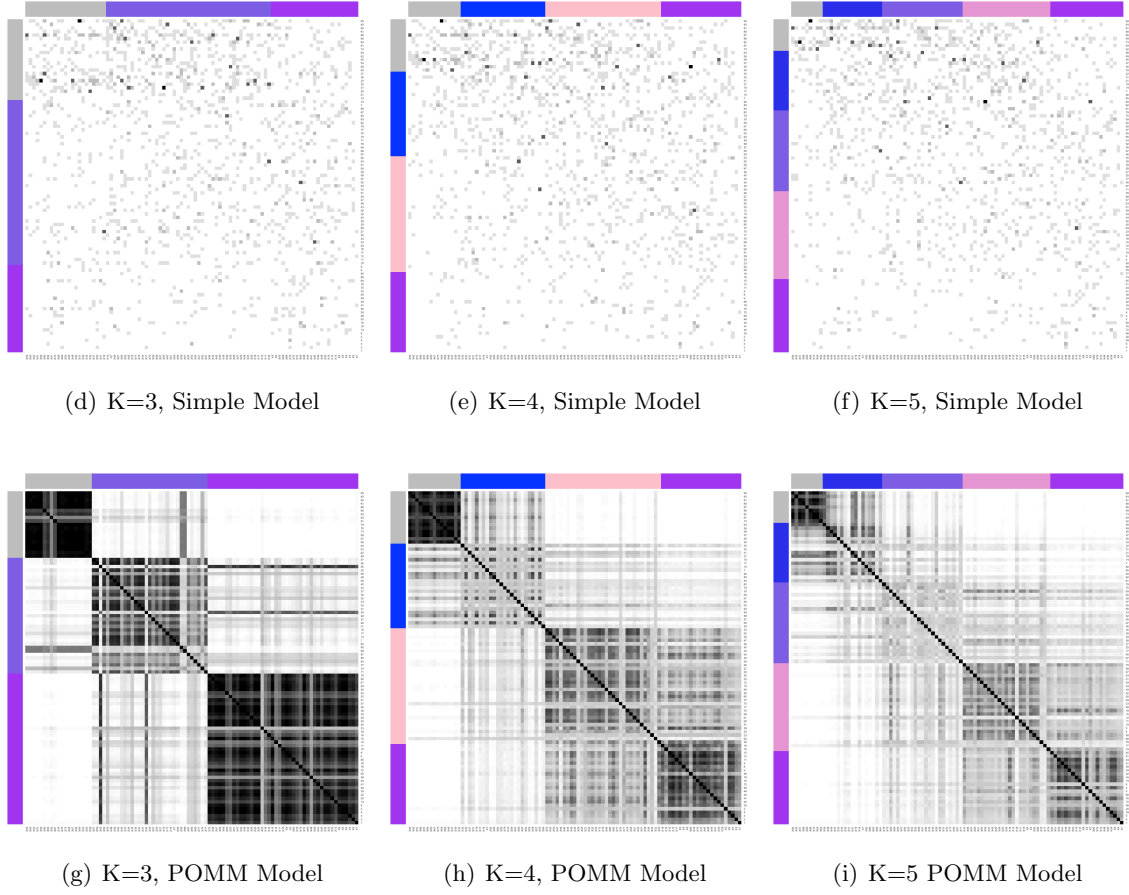


FIGURE 8. Adjacency Matrices simulated via the POMM Model

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

Method	$\hat{\theta}$			95% CI interval		
	(a)	(b)	(c)	(a)	(b)	(c)
S	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

11.2.1. *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 on 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.

12. APPENDIX I: ESTIMATION DETAILS

12.1. Updating \mathbf{z} . To update \mathbf{z} we propose a new label for each node, we evaluate the accept/reject move by computing the ratio r as follows:

$$(18) \quad r = \frac{\prod_{i < j} \binom{n_{ij}}{y_{ij}} p_{z'_i z'_j}^{y_{ij}} \cdot (1 - p_{z'_i z'_j})^{n_{ij} - y_{ij}} \cdot \frac{\Gamma(\gamma_0)\Gamma(n+1)}{\Gamma(n+\gamma_0)} \cdot \prod_{k=1}^K \frac{\Gamma(n'_k + \gamma_k)}{\Gamma(\gamma_k)\Gamma(n'_k + 1)}}{\prod_{i < j} \binom{n_{ij}}{y_{ij}} p_{z_i z_j}^{y_{ij}} \cdot (1 - p_{z_i z_j})^{n_{ij} - y_{ij}} \cdot \frac{\Gamma(\gamma_0)\Gamma(n+1)}{\Gamma(n+\gamma_0)} \cdot \prod_{k=1}^K \frac{\Gamma(n_k + \gamma_k)}{\Gamma(\gamma_k)\Gamma(n_k + 1)}}$$

$$(19) \quad = \frac{\prod_{i < j} p_{z'_i z'_j}^{y_{ij}} \cdot (1 - p_{z'_i z'_j})^{n_{ij} - y_{ij}} \cdot \prod_{k=1}^K \frac{\Gamma(n'_k + \gamma_k)}{\Gamma(\gamma_k)\Gamma(n'_k + 1)}}{\prod_{i < j} p_{z_i z_j}^{y_{ij}} \cdot (1 - p_{z_i z_j})^{n_{ij} - y_{ij}} \cdot \prod_{k=1}^K \frac{\Gamma(n_k + \gamma_k)}{\Gamma(\gamma_k)\Gamma(n_k + 1)}}$$

Passing to the log:

$$\begin{aligned} \log(r) &= \log \left(\prod_{i < j} p_{z'_i z'_j}^{y_{ij}} \cdot (1 - p_{z'_i z'_j})^{n_{ij} - y_{ij}} \cdot \prod_{k=1}^K \frac{\Gamma(n'_k + \gamma_k)}{\Gamma(\gamma_k)\Gamma(n'_k + 1)} \right) \\ &\quad - \log \left(\prod_{i < j} p_{z_i z_j}^{y_{ij}} \cdot (1 - p_{z_i z_j})^{n_{ij} - y_{ij}} \cdot \prod_{k=1}^K \frac{\Gamma(n_k + \gamma_k)}{\Gamma(\gamma_k)\Gamma(n_k + 1)} \right) \\ &= \sum_{i < j} \left(y_{ij} \cdot \log p_{z'_i z'_j} + (n_{ij} - y_{ij}) \cdot \log (1 - p_{z'_i z'_j}) \right) \\ &\quad + \sum_{k=1}^K \left(\log (\Gamma(n'_k + \gamma_k)) - \log (\Gamma(\gamma_k)) - \log (\Gamma(n'_k + 1)) \right) \\ &\quad - \sum_{i < j} \left(y_{ij} \cdot \log p_{z_i z_j} + (n_{ij} - y_{ij}) \cdot \log (1 - p_{z_i z_j}) \right) \\ &\quad - \sum_{k=1}^K \left(\log (\Gamma(n_k + \gamma_k)) - \log (\Gamma(\gamma_k)) - \log (\Gamma(n_k + 1)) \right) \end{aligned} \quad (20)$$

12.2. Updating \mathbf{P} . To update \mathbf{P} and α we propose a new label for each node, we evaluate the accept/reject move by computing the ratio r as follows:

$$(21) \quad r = \frac{\prod_{i < j} \binom{n_{ij}}{y_{ij}} p_{z_i z_j}'^{y_{ij}} \cdot (1 - p_{z_i z_j}')^{n_{ij} - y_{ij}} \cdot \prod_{k=1}^K \left(\frac{1}{y'^{(k+1)} - y'^{(k)}} \right)^{|L'^{(k)}|}}{\prod_{i < j} \binom{n_{ij}}{y_{ij}} p_{z_i z_j}^{y_{ij}} \cdot (1 - p_{z_i z_j})^{n_{ij} - y_{ij}} \cdot \prod_{k=1}^K \left(\frac{1}{y^{(k+1)} - y^{(k)}} \right)^{|L^{(k)}|}}$$

$$(22)$$

Passing to the log:

Algorithm 2 Updating z step

```

1: for  $i \leftarrow 1$  to  $N$  do
2:   Sample new_label from  $1, \dots, K$ 
3:   Set  $z' \leftarrow z$  with the  $i$ -th element replaced by new_label
4:   Compute new victory probabilities  $p_{z'_i z'_j}$  using  $z'$ 
5:   Compute probability ratio  $\log(r)$  using  $p_{z'_i z'_j}$  and  $p_{z_i z_j}$ 
6:   Set  $\alpha_r \leftarrow \min(1, r)$ 
7:   Sample  $u$  from a uniform distribution on  $(0, 1)$ 
8:   if  $u < \alpha_r$  then
9:     Update  $z$  to  $z'$ 
10:    Update  $p_{z_i z_j}$  to  $p_{z'_i z'_j}$ 
11:    Increment  $\text{acc.count}_z$ 
12:   end if
13:   Store  $z_{\text{current}}$  in  $z.\text{container}$ 
14: end for

```

(23)

$$\log(r) = \sum_{i < j} \left(y_{ij} \cdot \log p'_{z_i z_j} + (n_{ij} - y_{ij}) \cdot \log(1 - p'_{z_i z_j}) \right) - \sum_{k=1}^K |L'^{(k)}| \cdot \log \left(y'^{(k+1)} - y'^{(k)} \right)$$

(24)

$$- \sum_{i < j} \left(y_{ij} \cdot \log p_{z_i z_j} + (n_{ij} - y_{ij}) \cdot \log(1 - p_{z_i z_j}) \right) + \sum_{k=1}^K |L^{(k)}| \cdot \log \left(y^{(k+1)} - y^{(k)} \right)$$

13. APPENDIX II: POMM PRIOR CHECKS

13.1. Prior predictive check.

13.2. MLE check.

Algorithm 3 Updating P step

```

1:  $j \leftarrow 1$ 
2: while  $j \leq N_{iter}$  do
3:   Sample  $\alpha'$  from a truncated normal distribution
4:   Generate a new proposal matrix  $P'$ 
5:   Compute new victory probabilities  $p'_{z_i z_j}$  using  $P'$  and  $z_{current}$ 
6:   Compute probability ratio  $\log(r)$  using  $p'_{z_i z_j}$  and  $p_{z_i z_j}$ 
7:   Set  $\alpha_r \leftarrow \min(1, r)$ 
8:   Sample  $u$  from a uniform distribution on  $(0, 1)$ 
9:   if  $u < \alpha_r$  then
10:    Update  $\alpha$  to  $\alpha'$ 
11:    Update  $P$  to  $P'$ 
12:    Update  $p_{z_i z_j}$  to  $p'_{z_i z_j}$ 
13:    Increment  $acc.count_p$ 
14:   end if
15:   Store  $P$  in  $P.container$ 
16:   Store  $\alpha$  in  $\alpha.container$ 
17:    $j \leftarrow j + 1$ 
18: end while

```
