



## **Recurrent segmentation meets block models in temporal networks**

Chamalee Wickrama Arachchi, Nikolaj Tatti

University of Helsinki

**Discovery Science : 2022**

# Outline

1. Introduction and Background.
2. Problem Formulation.
3. Algorithms.
4. Monotonicity, SMAWK.
5. Experiments.

## Stochastic Block Model (SBM)

- SBM is a generative model for random graphs.
- It generates graphs containing communities such that connections between subsets of nodes are characterized with particular edge densities.
- Each node is assumed to belong to a particular group / cluster and the way nodes interact depends **only** on their respective cluster<sup>1</sup>.



**Figure:** : The above two graphs are the same graph re-organized and drawn from the SBM model with 1000 vertices, 5 balanced communities, within-cluster probability of .02 and across-cluster probability of .001<sup>2</sup>.

---

<sup>1</sup>Wilkinson D.J Lee C. A review of stochastic block models and extensions for graph clustering. Applied Network Science, 4(122), 2019.

<sup>2</sup>Emmanuel Abbe, Community detection and stochastic block models: recent developments, 2017.

## ■ Temporal Networks.

Unlike static networks, real-world networks are often timestamped.

Ex. Time stamped interactions between actors.

- ▶ X1 sends an SMS to Y1 at time  $t_1$ .
- ▶ X2 sends an email to Y1 at time  $t_2$ .
- ▶ X4 likes/answers to Y3's post at time  $t_3$ .
- ▶ X1 phones to Y4 at time  $t_4$ .
- ▶ and also: citations (patents, articles), web links, tweets, moving objects, etc.

## ■ Recurrent Activities.

Many temporal networks exhibit cyclic or repeating behaviour.

They might have *shared* parameters.

Ex. Genomic sequence might be generated with  $h$  hidden sources such that some of the sources contribute for multiple segments.

## Recurrent Segmentation Problem

### Given

Sequence of edges  
with timestamps.

$R$  = No of groups,  
 $K$  = No of segments,  
 $H$  = No of levels.

### Goal

Partition nodes using  
SBM core.

Segment the timeline into  
 $K$  segments.

Segments should have  $H$   
distinct levels.

### Model

Each node is assumed to  
belong to a particular  
group; extension of SBM  
to dynamic networks.

Edges are modelled with a  
Non-homogeneous  
Poisson process.

Timeline of activities has  
the support for recurrent  
activities.

Maximize the likelihood.

### Key questions

Is problem NP-hard?

Linear time algorithm?

## Toy Example

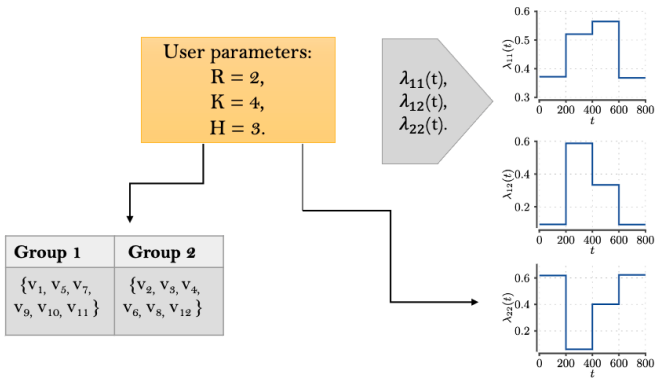


Figure: Toy Example

## Problem Formulation

- $R$  - Number of groups.
- $K$  - Number of segments.
- $H$  - Number of levels.
- $P_i$  - Set of vertices in  $i$ th group.
- $T_k$  -  $k$ th time segment.
- $\lambda_{ijk}$  - Poisson parameter between  $i$ th group and  $j$ th group in  $k$ th segment.
- $\Lambda = \{\lambda_{ijk}\}$ .

$$\ell(u, v, T, \lambda) = c(u, v, T) \log \lambda - \lambda \Delta(T)$$

$$\ell(U, W, T, \lambda) = \sum_{u, w \in U \times W} \ell(u, w, T, \lambda)$$

$$\ell(\mathcal{P}, \mathcal{T}, \Lambda) = \sum_{i=1}^R \sum_{j=i}^R \sum_{k=1}^K \ell(P_i, P_j, T_k, \lambda_{ijk})$$

## Algorithms

Searching for optimal groups and segments is computationally intractable. Therefore, we split the problem into three subproblems and alternately solve each subproblem.

- Optimize groups.
- Optimize Poisson parameters.
- Optimize segmentation.

---

### Algorithm 1: Main loop of the algorithm

---

```

1  $\mathcal{P} \leftarrow$  random groups;
2  $\Lambda \leftarrow$  random values;
3  $\mathcal{T}, g \leftarrow \text{FINDSEGMENTS}(\mathcal{P}, \Lambda);$ 
4  $\Lambda \leftarrow \text{UPDATELAMBDA}(\mathcal{P}, \mathcal{T}, g);$ 
5 while convergence do
6    $\mathcal{P} \leftarrow \text{FINDGROUPS}(\mathcal{P}, \Lambda, \mathcal{T}, g);$ 
7    $\Lambda \leftarrow \text{UPDATELAMBDA}(\mathcal{P}, \mathcal{T}, g);$ 
8    $\mathcal{T}, g \leftarrow \text{FINDSEGMENTS}(\mathcal{P}, \Lambda);$ 
9    $\Lambda \leftarrow \text{UPDATELAMBDA}(\mathcal{P}, \mathcal{T}, g);$ 

```

---



## Subproblem 1 : Optimize Groups

- Greedy algorithm.
- Based on likelihood gain, each node is assigned to a group one at a time while group assignments of other nodes are fixed.
- A naive implementation of computing the log-likelihood gain for a single node may require  $\Theta(m)$  steps, which would lead in  $\Theta(nm)$  time as we need to test every node.
- However, constant terms can be omitted when computing loglikelihood gain.
- For each node, we can consider only adjacent edges and then precompute edge count array  $c[j, k]$  ( $j$ th group in  $k$ th segment) in  $\mathcal{O}(m)$  time.
- The computational complexity for FINDGROUPS is  $\mathcal{O}(m + Rn + R^2H + K)$ .

## Subproblem 2 : Optimize Poisson Parameters

- The optimal parameters are equal to

$$\lambda_{ijk} = \frac{c(P_i, P_j, T)}{|P_i \times P_j| \Delta(T_k)} \quad .$$

- Updating the parameters requires  $\mathcal{O}(m + R^2 H + K)$  time.

## Subproblem 3 : Optimize Segmentation

Naive:

- Given a sequence of edges ordered by timestamp (Chronological order), use Dynamic Programming (DP) for segmentation.
- DP recurrence formula for sequence segmentation:

$$s[1, m; k] = \max_{k \leq j \leq m-1} s[1, j; k-1] + y(j+1, m; 1)$$

- Fill DP table ( $m \times k$ ) and store also maximizing  $j$  values such that we can recover the optimal segmentation by doing a trace back.
- Computational complexity is  $\mathcal{O}(m^2 K)$ .

Speedup:

- Use SMAWK to speed up the running time.
- Quadratic time to linear time using SMAWK.
- Due to totally monotonicity of matrix.

## Totally Monotone Property

- We call that  $x$  is *totally monotone*,  
if  $x(i_1, j_1) > x(i_1, j_2)$ , then  $x(i_2, j_1) > x(i_2, j_2)$  for any  $i_1 < i_2$  and  $j_1 < j_2$ .

		j1				j2	
i1		$x(i_1, j_1)$		>		$x(i_1, j_2)$	
				implies			
i2		$x(i_2, j_1)$		>		$x(i_2, j_2)$	

Figure: Totally monotone matrix

## SMAWK

- Due to totally - monotonicity.
- SMAWK can find the row maxima in every row in  $\mathcal{O}(m)$  time on  $n \times m$  wide totally monotone matrix.
- SMAWK recursively computes the leftmost row maxima in every other row, and then fill in the remaining row maxima in  $\mathcal{O}(m)$  time.

## Segmentation Algorithm

---

**Algorithm 3:** Algorithm FINDSEGMENTS( $\mathcal{P}, \Lambda$ ) for finding optimal segmentation for fixed groups  $\mathcal{P}$  and parameters  $\Lambda$

---

```

1   $t_{min} \leftarrow \min \{t \mid (u, v, t) \in E\};$ 
2   $f[e, h] \leftarrow$  log-likelihood of a segment  $[t_{min}, t(e)]$  using parameters  $\lambda_{..h}$ ;
3  foreach  $e \in E$  do  $o[e, 1] \leftarrow \max_h f[e, h]$  ;
4  foreach  $k = 2, \dots, K$  do
5       $x(s, e; h) \leftarrow o[s, k-1] + f[e, h] - f[s, h];$ 
6      foreach  $h = 1, \dots, H$  do
7           $z[e, h] \leftarrow \arg \max_s x(s, e; h)$  for each  $e \in E$  (use SMAWK);
8           $o[e, k] \leftarrow \max_h x(z[e, h], e; h)$  for each  $e \in E$ ;
9           $r[e, k] \leftarrow \arg \max_h x(z[e, h], e; h);$ 
10          $q[e, k] \leftarrow z[e, r[e, k]];$ 
11 Use  $r$  and  $q$  to recover the optimal segmentation  $(T_1, \dots, T_K)$  and the level mapping  $g$  ;
12 return  $(T_1, \dots, T_K), g;$ 

```

---

Figure: FINDSEGMENTS

## SMAWK to Speedup DP

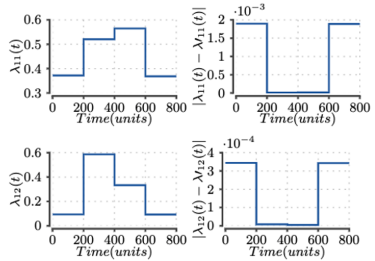
- $o[e, k]$  : log-likelihood of  $k$ -segmentation up to  $e$ .
- $y(s, e; h)$  : log-likelihood of a segment  $(t(s), t(e)]$  using the  $h$ th level of Poisson parameters.

$$o[e, k] = \max_h \max_s o[s, k-1] + y(s, e; h)$$

- The function  $x(s, e) = o[s, k-1] + y(s, e; h)$  is totally monotone; for each fixed  $h$ .
- To compute  $x$  in  $\mathcal{O}(1)$ , we precompute  $f[e, h]$  which is the log likelihood from epoch to  $e$  using corresponding  $h$  levels params;  $y(s, e; h) = f[e, h] - f[s, h]$ .
- SMAWK computes  $x$  by lazy fashion.
- Optimal segmentation is found after iterating over all  $h$  values.
- FINDSEGMENTS runs in  $\mathcal{O}(mKH + HR^2)$  time since we need to call SMAWK  $\mathcal{O}(HK)$  times.

## Experiments with Synthetic Datasets

Dataset	$n$	$m$	$R$	$K$	$H$	$LL_1$	$G$	$LL_2$	$I$	$CT$
Synthetic-1	50	76332	2	2	2	0.95	1	0.94	2	2.81s
Synthetic-2	30	95889	3	3	3	0.94	1	0.94	3	5.36s
Synthetic-3	20	65056	3	3	3	0.97	1	0.97	3	3.91s
Synthetic-4	60	537501	3	4	3	0.94	1	0.93	3	23.13s
Synthetic-5	10	33475	2	10	5	0.91	1	0.91	4	10.27s

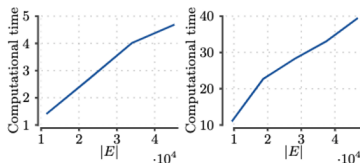


**Figure:** Table: Experiments for the synthetic datasets. Here,  $n$  is the number of nodes,  $m$  is the number of edges,  $R$  is the number of groups,  $K$  is the number of segments,  $H$  is the number of levels,  $LL_1$  is the initial normalized log-likelihood,  $G$  is the Rand index,  $LL_2$  is the discovered normalized log-likelihood,  $I$  is the number of iterations, and  $CT$  is the computational time., Plot: Right : True parameters  $\lambda_{11}(t)$ ,  $\lambda_{12}(t)$  (left) and their difference to the estimated parameters  $\lambda_{\hat{11}}(t)$ ,  $\lambda_{\hat{12}}(t)$  (right) for the *Synthetic-4* dataset.  $\lambda_{12}(t)$  implies the Poisson process parameter between group 1 and group 2 as a function of time. This experiment is done with  $n = 60$ ,  $m = 537501$ ,  $R = 3$ ,  $K = 4$ , and  $H = 3$ .



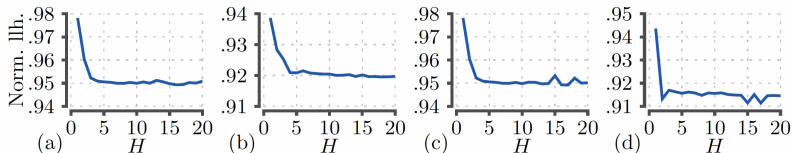
## Experiments with Real-world Datasets

<i>Dataset</i>	<i>n</i>	<i>m</i>	<i>R</i>	<i>K</i>	<i>H</i>	<i>LL</i>	<i>I</i>	<i>CT</i>
<i>Email-Eu-1</i>	309	61046	3	10	7	0.89	12	188s
<i>Email-Eu-2</i>	162	46772	4	8	7	0.87	9	177s
<i>MathOverflow</i>	21688	107581	2	3	2	0.91	20	263s
<i>CollegeMsg</i>	1899	59835	3	8	5	0.87	19	662s
<i>MOOC</i>	7047	411749	2	3	2	0.81	6	208s
<i>Bitcoin</i>	3783	24186	3	10	10	0.91	7	115s
<i>Santander</i>	735	33116	3	7	5	0.94	20	60s



**Figure:** Table: Sizes and computational times for the real-time datasets. Here,  $n$  is the number of nodes,  $m$  is the number of edges,  $R$  is the number of groups,  $K$  is the number of segments,  $H$  is the number of levels,  $LL$  is the normalized log-likelihood,  $I$  is the number of iterations, and  $CT$  is the computational time in seconds., Plot : Computational time (in seconds) as a function of number of temporal edges ( $|E|$ ) for the *Synthetic* datasets (left) and *Email-Eu-2* dataset (right). This experiment is done for  $R = 3$ ,  $K = 5$ , and  $H = 3$ .

## Likelihood vs Number of Levels ( $H$ )



**Figure:** Normalized log-likelihood as a function of number of levels ( $H$ ) for the *Santander* dataset (a), *bitcoin* dataset (b), *Synthetic-5* dataset (c), and *Email-Eu-1* dataset (d). This experiment is done for  $R = 2$ ,  $K = 20$ , and  $H = 1, \dots, 20$ .

**Any Questions???**

**Thank you!!!**