

Generative Model for Change Point Detection in Dynamic Graphs

Anonymous authors
Paper under double-blind review

Abstract

1 This paper proposes a generative model to detect change points in time series of graphs.
2 The proposed framework consists of learnable prior distributions for low-dimensional graph
3 representations and of a decoder that can generate graphs from the latent representations.
4 The prior distributions of the latent spaces are learned from the observed data as empirical
5 Bayes and generative model is employed to assist multiple change point detection. Specifi-
6 cally, the model parameters are learned via maximum approximate likelihood, with a Group
7 Fused Lasso regularization on the prior parameters. The optimization problem is then solved
8 via Alternating Direction Method of Multipliers and Langevin Dynamics are recruited for
9 posterior inference. Simulation studies show good performance of the generative model in
10 supporting change point detection, and real data experiments yield change points that align
11 with significant events.

12 1 Introduction

13 Networks are often used to represent relational phenomena in numerous domains (Dwivedi et al., 2021; He
14 et al., 2023; Han et al., 2023) and relational phenomena by nature progress in time. In recent decades, a
15 plethora of network models has been proposed to analyze the interaction between objects or people over
16 time, including Temporal Exponential-family Random Graph Model (Hanneke et al., 2010; Krivitsky &
17 Handcock, 2014), Stochastic Actor-Oriented Model (Snijders, 2001; Snijders et al., 2010), and Relational
18 Event Model (Butts, 2008; Butts et al., 2023). Although these models incorporate the temporal aspect for
19 network analysis, network evolution is usually time-heterogeneous. Without taking the structural changes
20 across dynamic networks into consideration, learning from the time series may lead to ambiguity. Hence, it
21 is practical for researchers to localize the change points before studying the evolving networks.

22 More recently, various methodologies have been proposed to detect change points in dynamic networks.
23 Chen et al. (2020) and Shen et al. (2023) employed embedding methods to detect both anomalous graphs
24 and vertices in time series of networks. Park & Sohn (2020) combined the multi-linear tensor regression
25 model with a hidden Markov model, detecting changes based on the transition between the hidden states.
26 Sulem et al. (2023) learned a graph similarity function using a Siamese graph neural network to differentiate
27 the graphs before and after a change point. Zhao et al. (2019) developed a screening algorithm that is
28 based on an initial graphon estimation to detect change points. Huang et al. (2020) utilized the singular
29 values of the Laplacian matrices as graph embedding to detect the differences across time. Chen & Zhang
30 (2015), Chu & Chen (2019), and Song & Chen (2022a) proposed a non-parametric approach to delineate
31 the distributional differences over time, and Garreau & Arlot (2018) and Song & Chen (2022b) exploited
32 the patterns in high dimensions via a kernel-based method. Zhang et al. (2024) jointly trained a Variational
33 Graph Auto-Encoder and Gaussian Mixture Model to detect the change points.

34 Inherently, network structures are complex due to highly dyadic dependency. Acquiring a low dimensional
35 representation of a graph can summarize the enormous amount of individual relations to promote downstream
36 analysis (Gallagher et al., 2021). In particular, Sharifnia & Saghaei (2022) and Kei et al. (2023) proposed to
37 detect the structural changes using an Exponential-family Random Graph Model. Yet they relied on user-
38 specified network statistics, which are usually not known to the modeler a priori. Moreover, Larroca et al.

(2021), Marengo et al. (2022), and Gong et al. (2023) developed different latent space models for dynamic graphs to detect changes, but they focused on node level representation, which may not be powerful enough to capture the information of the entire graph. Consequently, we aim to infer the graph level representations that induce the structural changes to facilitate the detection.

On the other hand, generative models recently showed promising results in myriad applications, such as text generation with Large Language Model (Devlin et al., 2018; Lewis et al., 2019) and image generation with Diffusion Model (Ho et al., 2020; Rombach et al., 2022). Similarly, we aim to explore how generative models can assist change point detection in dynamic graphs. In particular, Simonovsky & Komodakis (2018) proposed a Graph Variational Auto-Encoder (VAE) for graph generation, with a zero-mean Gaussian prior to regularize the latent space of the graph level representation. In the VAE framework (Kingma, 2013; Kipf & Welling, 2016), regularization via Kullback Leibler divergence arises from the Evidence Lower Bound for the marginal likelihood, encouraging the approximate posterior to be close to the zero-mean Gaussian prior. In contrast, we focus on learning the mean of the Gaussian prior at each time point and we apply Group Fused Lasso regularization to promote sparsity in the sequential differences of the prior parameters, effectively smoothing out minor fluctuations and highlighting significant change points.

To tackle the challenges in dynamic graphs and to employ recent advances in generative models, we make the following contributions in this manuscript:

- We learn graph level representations of network structures to facilitate change point detection. We impose a prior distribution to the representation at each time point and a multivariate total variation regularization to the sequential differences of the prior parameters. The prior distributions and a graph decoder are jointly learned via maximum approximate likelihood.
- We derive an Alternating Direction Method of Multipliers (ADMM) procedure to solve the optimization problem. Without using an encoder, the prior distributions and the graph decoder are learned by inferring from the posterior distribution via Langevin Dynamics. Experiments show good performance of the generative model in supporting change point detection.

The rest of the manuscript is organized as follows. Section 2 specifies the proposed framework. Section 3 presents the objective function with Group Fused Lasso regularization and the ADMM procedure to solve the optimization problem. Section 4 discusses change points localization and model selection. Section 5 illustrates the proposed method on simulated and real data. Section 6 concludes the work with a discussion and potential future developments.

2 Generative Model for Change Point Detection

2.1 Model Specification

For a node set $N = \{1, 2, \dots, n\}$, we use an adjacency matrix $\mathbf{y} \in \{0, 1\}^{n \times n}$ to represent a graph. We denote the set of all possible node pairs as $\mathbb{Y} = N \times N$. In the adjacency matrix, $\mathbf{y}_{ij} = 1$ indicates an edge between nodes i and j , while $\mathbf{y}_{ij} = 0$ indicates no edge. The relations can be either directed or undirected. The undirected variant has $\mathbf{y}_{ij} = \mathbf{y}_{ji}$ for all $(i, j) \in \mathbb{Y}$.

Denote \mathbf{y}^t as a network at a discrete time point t . The observed data is a sequence of networks $\mathbf{y}^1, \dots, \mathbf{y}^T$. For each network \mathbf{y}^t , we assume there is a latent variable $\mathbf{z}^t \in \mathbb{R}^d$ such that the network \mathbf{y}^t can be generated from the latent variable with the following graph decoder:

$$\mathbf{y}^t \sim P(\mathbf{y}^t | \mathbf{z}^t) = \prod_{(i,j) \in \mathbb{Y}} \text{Bernoulli}(\mathbf{y}_{ij}^t; \mathbf{r}_{ij}(\mathbf{z}^t))$$

where $\mathbf{r}_{ij}(\mathbf{z}^t) = P(\mathbf{y}_{ij}^t = 1 | \mathbf{z}^t)$ is the Bernoulli parameter for dyad (i, j) and it is elaborated in Section 2.2. Conditioning on the latent variable \mathbf{z}^t , we assume the network \mathbf{y}^t is dyadic independent.

We also impose a learnable prior distribution to the latent variable as

$$\mathbf{z}^t \sim P(\mathbf{z}^t) = \mathcal{N}(\mathbf{z}^t; \boldsymbol{\mu}^t, \mathbf{I}_d)$$

77 where $\boldsymbol{\mu}^t \in \mathbb{R}^d$ and \mathbf{I}_d is an identity matrix. With the graph decoder $P(\mathbf{y}^t|\mathbf{z}^t)$, we consider $\mathbf{z}^t \in \mathbb{R}^d$ as a
 78 graph level representation for $\mathbf{y}^t \in \{0, 1\}^{n \times n}$.

79 2.2 Graph Decoder

The graph decoder $P(\mathbf{y}^t|\mathbf{z}^t)$ is formulated with a Bernoulli parameter for dyad (i, j) as

$$\mathbf{r}_{ij}(\mathbf{z}^t) = P(\mathbf{y}_{ij}^t = 1|\mathbf{z}^t) = \mathbf{g}_{ij}(\mathbf{h}(\mathbf{z}^t)).$$

80 The function $\mathbf{h}(\cdot)$ is parameterized by neural networks with $\mathbf{h} : \mathbb{R}^d \rightarrow \mathbb{R}^{n \times n}$. The function $\mathbf{g}(\cdot)$ is the
 81 element-wise sigmoid function with $\mathbf{g} : \mathbb{R}^{n \times n} \rightarrow [0, 1]^{n \times n}$.

In particular, we use multi-layer perceptrons, transferring the latent variable $\mathbf{z}^t \in \mathbb{R}^d$ to $\mathbf{U}^t \in \mathbb{R}^{n \times k}$ and
 $\mathbf{V}^t \in \mathbb{R}^{n \times k}$, respectively. We let the latent dimension d and k be smaller than the number of nodes n , and

$$\mathbf{h}(\mathbf{z}^t) = \begin{cases} \mathbf{U}^t \mathbf{V}^{t\top} \in \mathbb{R}^{n \times n}, & \text{for directed network,} \\ \mathbf{U}^t \mathbf{U}^{t\top} \in \mathbb{R}^{n \times n}, & \text{for undirected network.} \end{cases}$$

82 The graph decoder via matrix multiplication is common in the literature (Kipf & Welling, 2016; Hamilton
 83 et al., 2017; Pan et al., 2018). Comparing to a decoder that directly outputs the graph $\mathbf{y}^t \in \{0, 1\}^{n \times n}$, the
 84 decoder via matrix multiplication reduces the number of parameters in the neural networks and helps avoid
 85 over-parameterization, which is crucial for graphs that are sparse.

86 Figure 1 gives an overview of the proposed framework. Implicitly, the graph level representation \mathbf{z}^t progresses
 87 to node level representations \mathbf{U}^t and \mathbf{V}^t as an intermediate step, before the generation of network \mathbf{y}^t . The
 88 graph decoder $P_\phi(\mathbf{y}^t|\mathbf{z}^t)$ with neural network parameter ϕ is shared across $t = 1, \dots, T$. It is worth pointing
 89 out the simplicity of our framework, without the need of encoders.

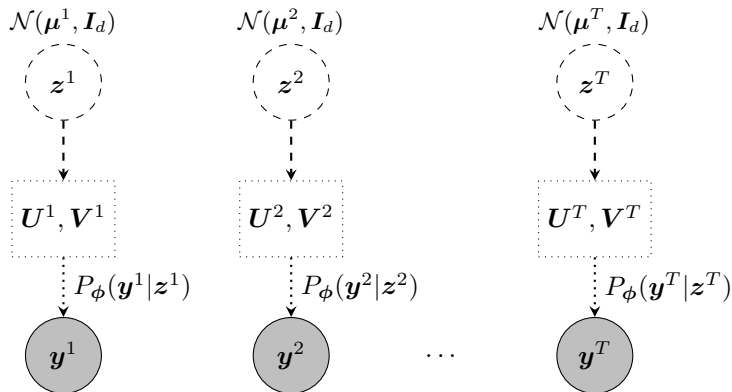


Figure 1: An overview of prior distributions and graph decoder.

90 2.3 Change Points

Anchored on the proposed framework, we can now specify the change points to be detected, in terms of the
 prior parameters $\boldsymbol{\mu}^t \in \mathbb{R}^d$ for $t = 1, \dots, T$. Let $\{C_k\}_{k=0}^{K+1} \subset \{1, 2, \dots, T\}$ be a collection of ordered change
 points with $1 = C_0 < C_1 < \dots < C_K < C_{K+1} = T$ such that

$$\begin{aligned} \boldsymbol{\mu}^{C_k} &= \boldsymbol{\mu}^{C_{k+1}} = \dots = \boldsymbol{\mu}^{C_{k+1}-1}, \quad k = 0, \dots, K, \\ \boldsymbol{\mu}^{C_k} &\neq \boldsymbol{\mu}^{C_{k+1}}, \quad k = 0, \dots, K-1, \quad \text{and} \quad \boldsymbol{\mu}^{C_{K+1}} = \boldsymbol{\mu}^{C_K}. \end{aligned}$$

91 The associated multiple change point detection problem comprises recovering the collection $\{C_k\}_{k=1}^K$ from a
 92 sequence of observed networks $\{\mathbf{y}^t\}_{t=1}^T$, where the number of change points K is also unknown.

93 To facilitate change point detection for $\{\mathbf{y}^t\}_{t=1}^T$ in the data space, we turn to learn the prior parameters
 94 $\{\boldsymbol{\mu}^t\}_{t=1}^T$ in the latent space. Intuitively, the consecutive prior parameters are similar when no change occurs,
 95 but they are different when a change emerges. For notational simplicity, we denote $\boldsymbol{\mu} \in \mathbb{R}^{T \times d}$ as a matrix
 96 where the t -th row corresponds to $\boldsymbol{\mu}^t \in \mathbb{R}^d$ with $t = 1, \dots, T$.

97 3 Learning and Inference

98 3.1 Learning Priors from Dynamic Graphs

99 Inspired by Vert & Bleakley (2010) and Bleakley & Vert (2011), we formulate the change point detection
 100 problem as a Group Fused Lasso problem (Alaíz et al., 2013). Denote the log-likelihood of the distribution
 101 for $\mathbf{y}^1, \dots, \mathbf{y}^T$ as $l(\boldsymbol{\phi}, \boldsymbol{\mu})$. We want to solve

$$\hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\mu}} = \arg \min_{\boldsymbol{\phi}, \boldsymbol{\mu}} -l(\boldsymbol{\phi}, \boldsymbol{\mu}) + \lambda \sum_{t=1}^{T-1} \|\boldsymbol{\mu}^{t+1} - \boldsymbol{\mu}^t\|_2 \quad (1)$$

102 where $\lambda > 0$ is a tuning parameter for the Group Fused Lasso penalty term.

103 The Group Fused Lasso penalty is useful for change point detection because it enforces piecewise constant
 104 patterns in the learned parameters by minimizing the total variation. Specifically, the regularization term,
 105 expressed as the sum of the ℓ_2 norms, encourages sparsity of the differences $\boldsymbol{\mu}^{t+1} - \boldsymbol{\mu}^t \in \mathbb{R}^d$, while allowing
 106 multiple coordinates across the d dimensional differences to change at the same time t . The latter is often
 107 referred as a grouping effect that could not be achieved with the ℓ_1 penalty of the differences. Furthermore,
 108 since the regularization is imposed on the prior parameters that relate to the likelihood of the data, the
 109 learned priors incorporate the structural changes from the observed graphs into the latent space. In summary,
 110 by penalizing the sum of sequential differences, the proposed framework focuses on capturing meaningful
 111 structural changes and smoothing out minor variations.

112 To solve the optimization problem in (1) that involves latent variables, we need to manipulate the objective
 113 function accordingly. We first introduce a slack variable $\boldsymbol{\nu} \in \mathbb{R}^{T \times d}$ where $\boldsymbol{\nu}^t \in \mathbb{R}^d$ denotes the t -th row of
 114 matrix $\boldsymbol{\nu}$, and we can rewrite the original problem as a constrained optimization problem:

$$\begin{aligned} \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\mu}} = \arg \min_{\boldsymbol{\phi}, \boldsymbol{\mu}} & -l(\boldsymbol{\phi}, \boldsymbol{\mu}) + \lambda \sum_{t=1}^{T-1} \|\boldsymbol{\nu}^{t+1} - \boldsymbol{\nu}^t\|_2 \\ & \text{subject to } \boldsymbol{\mu} = \boldsymbol{\nu}. \end{aligned} \quad (2)$$

Then the augmented Lagrangian can be defined as

$$\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{\rho}) = -l(\boldsymbol{\phi}, \boldsymbol{\mu}) + \lambda \sum_{t=1}^{T-1} \|\boldsymbol{\nu}^{t+1} - \boldsymbol{\nu}^t\|_2 + \text{tr}[\boldsymbol{\rho}^\top (\boldsymbol{\mu} - \boldsymbol{\nu})] + \frac{\kappa}{2} \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_F^2$$

115 where $\boldsymbol{\rho} \in \mathbb{R}^{T \times d}$ is the Lagrange multipliers and $\kappa > 0$ is the penalty parameter for the augmentation term.
 116 Let $\mathbf{w} = \kappa^{-1} \boldsymbol{\rho} \in \mathbb{R}^{T \times d}$ be the scaled dual variable, the augmented Lagrangian can be updated to

$$\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\nu}, \mathbf{w}) = -l(\boldsymbol{\phi}, \boldsymbol{\mu}) + \lambda \sum_{t=1}^{T-1} \|\boldsymbol{\nu}^{t+1} - \boldsymbol{\nu}^t\|_2 + \frac{\kappa}{2} \|\boldsymbol{\mu} - \boldsymbol{\nu} + \mathbf{w}\|_F^2 - \frac{\kappa}{2} \|\mathbf{w}\|_F^2. \quad (3)$$

In practice, gradient descent may not work well for an objective function with Group Fused Lasso penalty.
 We further introduce two variables $(\boldsymbol{\gamma}, \boldsymbol{\beta}) \in \mathbb{R}^{1 \times d} \times \mathbb{R}^{(T-1) \times d}$ to ease the optimization, by converting it into
 a Group Lasso problem (Yuan & Lin, 2006). They are defined as

$$\boldsymbol{\gamma} = \boldsymbol{\nu}^1 \quad \text{and} \quad \boldsymbol{\beta}_{t,\cdot} = \boldsymbol{\nu}^{t+1} - \boldsymbol{\nu}^t \quad \forall t = 1, \dots, T-1.$$

Reversely, the slack variable $\boldsymbol{\nu} \in \mathbb{R}^{T \times d}$ can be reconstructed as $\boldsymbol{\nu} = \mathbf{1}_{T,1}\boldsymbol{\gamma} + \mathbf{X}\boldsymbol{\beta}$, where \mathbf{X} is a $T \times (T-1)$ design matrix with $\mathbf{X}_{ij} = 1$ for $i > j$ and 0 otherwise. Substituting the $\boldsymbol{\nu}$ in (3) with $(\boldsymbol{\gamma}, \boldsymbol{\beta})$, we have

$$\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\gamma}, \boldsymbol{\beta}, \mathbf{w}) = -l(\boldsymbol{\phi}, \boldsymbol{\mu}) + \lambda \sum_{t=1}^{T-1} \|\boldsymbol{\beta}_{t,\cdot}\|_2 + \frac{\kappa}{2} \|\boldsymbol{\mu} - \mathbf{1}_{T,1}\boldsymbol{\gamma} - \mathbf{X}\boldsymbol{\beta} + \mathbf{w}\|_F^2 - \frac{\kappa}{2} \|\mathbf{w}\|_F^2.$$

Thus, we can derive the following Alternating Direction Method of Multipliers (ADMM) procedure (Boyd et al., 2011; Wang et al., 2019) to solve the constrained optimization problem in (2):

$$\boldsymbol{\phi}_{(a+1)}, \boldsymbol{\mu}_{(a+1)} = \arg \min_{\boldsymbol{\phi}, \boldsymbol{\mu}} -l(\boldsymbol{\phi}, \boldsymbol{\mu}) + \frac{\kappa}{2} \|\boldsymbol{\mu} - \boldsymbol{\nu}_{(a)} + \mathbf{w}_{(a)}\|_F^2, \quad (4)$$

$$\boldsymbol{\gamma}_{(a+1)}, \boldsymbol{\beta}_{(a+1)} = \arg \min_{\boldsymbol{\gamma}, \boldsymbol{\beta}} \lambda \sum_{t=1}^{T-1} \|\boldsymbol{\beta}_{t,\cdot}\|_2 + \frac{\kappa}{2} \|\boldsymbol{\mu}_{(a+1)} - \mathbf{1}_{T,1}\boldsymbol{\gamma} - \mathbf{X}\boldsymbol{\beta} + \mathbf{w}_{(a)}\|_F^2, \quad (5)$$

$$\mathbf{w}_{(a+1)} = \boldsymbol{\mu}_{(a+1)} - \boldsymbol{\nu}_{(a+1)} + \mathbf{w}_{(a)}, \quad (6)$$

117 where subscript a denotes the current ADMM iteration. We recursively implement the three updates until
 118 a convergence criterion is satisfied. Throughout the paper, details about the implementation are provided
 119 in Appendix 7.4.

120 3.2 Parameters Update

121 3.2.1 Updating $\boldsymbol{\mu}$ and $\boldsymbol{\phi}$

122 In this section, we derive the updates for the prior and graph decoder parameters. Denote the objective
 123 function in (4) as $\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\mu})$. Setting the gradients of $\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\mu})$ with respect to the prior parameter $\boldsymbol{\mu}^t \in \mathbb{R}^d$ to
 124 zeros, we have the following:

125 **Proposition 1.** *The solution for $\boldsymbol{\mu}^t$ at an iteration of our proposed ADMM algorithm is a weighted sum:*

$$\boldsymbol{\mu}^t = \frac{1}{1+\kappa} \mathbb{E}_{P(\mathbf{z}^t|\mathbf{y}^t)}(\mathbf{z}^t) + \frac{\kappa}{1+\kappa} (\boldsymbol{\nu}^t - \mathbf{w}^t) \quad (7)$$

126 *between the conditional expectation of the latent variable under the posterior distribution $P(\mathbf{z}^t|\mathbf{y}^t)$ and the*
 127 *difference between the slack and the scaled dual variables. The term $\mathbf{w}^t \in \mathbb{R}^d$ denotes the t -th row of the*
 128 *scaled dual variable $\mathbf{w} \in \mathbb{R}^{T \times d}$. The derivation is provided in Appendix 7.1.*

129 Moreover, the gradient of $\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\mu})$ with respect to the graph decoder parameter $\boldsymbol{\phi}$ is calculated as

$$\nabla_{\boldsymbol{\phi}} \mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\mu}) = - \sum_{t=1}^T \mathbb{E}_{P(\mathbf{z}^t|\mathbf{y}^t)} \left(\nabla_{\boldsymbol{\phi}} \log P(\mathbf{y}^t|\mathbf{z}^t) \right). \quad (8)$$

130 The parameter $\boldsymbol{\phi}$ can be updated efficiently through back-propagation.

Notably, calculating the solution in (7) and the gradient in (8) requires evaluating the conditional expectation under the posterior distribution $P(\mathbf{z}^t|\mathbf{y}^t) \propto P(\mathbf{y}^t|\mathbf{z}^t) \times P(\mathbf{z}^t)$. We employ Langevin Dynamics to sample from the posterior distribution, approximating the conditional expectations (Xie et al., 2017; 2018; Nijkamp et al., 2020; Pang et al., 2020). In particular, let subscript τ be the time step of the Langevin Dynamics and let δ be a small step size. Moving toward the gradient of the posterior with respect to the latent variable, the Langevin Dynamics to draw samples from the posterior distribution is achieved by iterating:

$$\begin{aligned} \mathbf{z}_{\tau+1}^t &= \mathbf{z}_{\tau}^t + \delta [\nabla_{\mathbf{z}^t} \log P(\mathbf{z}^t|\mathbf{y}^t)] + \sqrt{2\delta} \boldsymbol{\epsilon} \\ &= \mathbf{z}_{\tau}^t + \delta [\nabla_{\mathbf{z}^t} \log P_{\boldsymbol{\phi}}(\mathbf{y}^t|\mathbf{z}^t) - (\mathbf{z}_{\tau}^t - \boldsymbol{\mu}^t)] + \sqrt{2\delta} \boldsymbol{\epsilon} \end{aligned} \quad (9)$$

131 where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ is a random perturbation to the process. The derivation is provided in Appendix 7.2.

132 3.2.2 Updating γ and β

133 In this section, we derive the update in (5), which is equivalent to solving a Group Lasso problem. In
 134 particular, we decompose the slack variable ν to work with γ and β . With ADMM, the updates on γ and
 135 β do not require the observed network data $\{\mathbf{y}^t\}_{t=1}^T$.

136 By adapting the derivation in Bleakley & Vert (2011), we have the following for our proposed ADMM:

137 **Proposition 2.** [Bleakley & Vert, 2011] *The Group Lasso problem to update $\beta \in \mathbb{R}^{(T-1) \times d}$ is solved in
 138 a block coordinate descent manner, by iteratively applying the following equation to each row t :*

$$\beta_{t,\cdot} \leftarrow \frac{1}{\kappa \mathbf{X}_{\cdot,t}^\top \mathbf{X}_{\cdot,t}} \left(1 - \frac{\lambda}{\|\mathbf{b}_t\|_2} \right)_+ \mathbf{b}_t \quad (10)$$

where $(\cdot)_+ = \max(\cdot, 0)$ and

$$\mathbf{b}_t = \kappa \mathbf{X}_{\cdot,t}^\top (\boldsymbol{\mu}_{(a+1)} + \mathbf{w}_{(a)} - \mathbf{1}_{T,1} \gamma - \mathbf{X}_{\cdot,-t} \boldsymbol{\beta}_{-t,\cdot}).$$

139 The derivation is provided in Appendix 7.3.

The convergence of the procedure can be monitored by the Karush-Kuhn-Tucker conditions: for all $\beta_{t,\cdot} \neq \mathbf{0}$,

$$\lambda \frac{\beta_{t,\cdot}}{\|\beta_{t,\cdot}\|_2} - \kappa \mathbf{X}_{\cdot,t}^\top (\boldsymbol{\mu}_{(a+1)} + \mathbf{w}_{(a)} - \mathbf{1}_{T,1} \gamma - \mathbf{X} \boldsymbol{\beta}) = \mathbf{0},$$

and for all $\beta_{t,\cdot} = \mathbf{0}$,

$$\|-\kappa \mathbf{X}_{\cdot,t}^\top (\boldsymbol{\mu}_{(a+1)} + \mathbf{w}_{(a)} - \mathbf{1}_{T,1} \gamma - \mathbf{X} \boldsymbol{\beta})\|_2 \leq \lambda.$$

Lastly, for any $\boldsymbol{\beta}$, the minimum in $\gamma \in \mathbb{R}^{1 \times d}$ is achieved at

$$\gamma = (1/T) \mathbf{1}_{1,T} \cdot (\boldsymbol{\mu}_{(a+1)} + \mathbf{w}_{(a)} - \mathbf{X} \boldsymbol{\beta}).$$

140 In summary, the procedure to solve the problem in (2) via ADMM is presented in Algorithm 1. The steps to
 141 transform between ν and (γ, β) within an ADMM iteration are omitted for succinctness. The complexity of
 142 the proposed algorithm is at least of order $O(A(Tsl + BT + D(T-1)))$ with additional gradient calculation
 143 for neural networks in the sub-routines. Specifically, for each of the A iterations of ADMM, we update the
 144 prior parameters $\boldsymbol{\mu}^t$ for all T time points, and each update involves l steps of MCMC for s samples. Then we
 145 calculate the gradients for neural networks over the T time points and run B iterations of Adam optimizer.
 146 Lastly, we run D iterations of block coordinate descent for the $T-1$ sequential differences.

147 4 Change Point Localization and Model Selection

148 4.1 Change Point Localization

In this section, we provide two effective methods to localize the change points after parameter learning, and
 they can be used for different purposes. For the first approach, we can resort to the prior distribution where
 $\mathbf{z}^t \sim \mathcal{N}(\boldsymbol{\mu}^t, \mathbf{I}_d)$. When no change occurs or $\boldsymbol{\mu}^t - \boldsymbol{\mu}^{t-1} = \mathbf{0}$, we have $\mathbf{z}^t - \mathbf{z}^{t-1} \sim \mathcal{N}(\mathbf{0}, 2\mathbf{I}_d)$ and

$$u^t := \frac{1}{2} (\mathbf{z}^t - \mathbf{z}^{t-1})^\top (\mathbf{z}^t - \mathbf{z}^{t-1}) \sim \chi_d^2.$$

Furthermore, the mean of u^t over m samples follows a Gamma distribution:

$$\bar{u}^t \sim \Gamma(\theta = \frac{2}{m}, \xi = \frac{md}{2})$$

149 where θ and ξ are the respective scale and shape parameters.

Algorithm 1 Latent Space Group Fused Lasso

```

1: Input: learning iterations  $A, B, D$ , tuning parameter  $\lambda$ , penalty parameter  $\kappa$ , learning rates  $\eta$ , observed
   data  $\{\mathbf{y}^t\}_{t=1}^T$ , initialization  $\{\phi_{(1)}, \boldsymbol{\mu}_{(1)}, \boldsymbol{\gamma}_{(1)}, \boldsymbol{\beta}_{(1)}, \mathbf{w}_{(1)}\}$ 
2: for  $a = 1, \dots, A$  do
3:   for  $t = 1, \dots, T$  do
4:     draw  $s$  samples  $\mathbf{z}_1^t, \dots, \mathbf{z}_s^t$  from  $P(\mathbf{z}^t | \mathbf{y}^t)$  according to (9)
5:      $\boldsymbol{\mu}_{(a+1)}^t = \frac{1}{1+\kappa} (s^{-1} \sum_{i=1}^s \mathbf{z}_i^t) + \frac{\kappa}{1+\kappa} (\boldsymbol{\nu}^t - \mathbf{w}^t)$ 
6:   end for
7:   for  $b = 1, \dots, B$  do
8:      $\phi_{(b+1)} = \phi_{(b)} - \eta \times \nabla_{\phi} \mathcal{L}(\phi, \boldsymbol{\mu})$ 
9:   end for
10:  Set  $\tilde{\boldsymbol{\gamma}}^{(1)} = \boldsymbol{\gamma}_{(a)}$  and  $\tilde{\boldsymbol{\beta}}^{(1)} = \boldsymbol{\beta}_{(a)}$ 
11:  for  $d = 1, \dots, D$  do
12:    for  $t = 1, \dots, T - 1$  do
13:      Let  $\tilde{\boldsymbol{\beta}}_{t,\cdot}^{(d+1)}$  be updated according to (10)
14:    end for
15:     $\tilde{\boldsymbol{\gamma}}^{(d+1)} = (1/T) \mathbf{1}_{1,T} \cdot (\boldsymbol{\mu}_{(a+1)} + \mathbf{w}_{(a)} - \mathbf{X} \tilde{\boldsymbol{\beta}}^{(d+1)})$ 
16:  end for
17:  Set  $\boldsymbol{\gamma}_{(a+1)} = \tilde{\boldsymbol{\gamma}}^{(d+1)}$  and  $\boldsymbol{\beta}_{(a+1)} = \tilde{\boldsymbol{\beta}}^{(d+1)}$ 
18:   $\mathbf{w}_{(a+1)} = \boldsymbol{\mu}_{(a+1)} - \boldsymbol{\nu}_{(a+1)} + \mathbf{w}_{(a)}$ 
19: end for
20:  $\hat{\boldsymbol{\mu}} \leftarrow \boldsymbol{\mu}_{(a+1)}$ 
21: Output: learned prior parameters  $\hat{\boldsymbol{\mu}}$ 

```

As we capture the structural changes in the latent space, we can draw samples from the learned priors to reflect the sequential changes. In particular, for a time point t , we sample $\hat{\mathbf{z}}^t - \hat{\mathbf{z}}^{t-1}$ from $\mathcal{N}(\hat{\boldsymbol{\mu}}^t - \hat{\boldsymbol{\mu}}^{t-1}, 2\mathbf{I}_d)$, and we perform the same transformation:

$$v^t := \frac{1}{2} (\hat{\mathbf{z}}^t - \hat{\mathbf{z}}^{t-1})^\top (\hat{\mathbf{z}}^t - \hat{\mathbf{z}}^{t-1}).$$

150 Then we compare the mean of v^t over m samples with a quantile:

$$P(\bar{v}^t > q_{\text{thr}}) = 1 - \frac{\alpha}{T-1} \quad (11)$$

151 where q_{thr} is the $1 - \alpha/(T-1)$ quantile of the Gamma distribution for \bar{u}^t when no change occurs. We
 152 consider the time point t with $\bar{v}^t > q_{\text{thr}}$ as the detected change point.

For the second approach, we can utilize the localizing method from Kei et al. (2023), which is more robust in practice, as compared in the simulation study of Section 5.1. First, we calculate the differences between consecutive time points in $\hat{\boldsymbol{\mu}} \in \mathbb{R}^{T \times d}$ as

$$\Delta \hat{\boldsymbol{\mu}}^t = \|\boldsymbol{\mu}^t - \boldsymbol{\mu}^{t-1}\|_2 \quad \forall t \in [2, T].$$

153 Then we standardize the differences as

$$\Delta \hat{\boldsymbol{\zeta}}^t = \frac{\Delta \hat{\boldsymbol{\mu}}^t - \text{median}(\Delta \hat{\boldsymbol{\mu}})}{\text{std}(\Delta \hat{\boldsymbol{\mu}})} \quad \forall t \in [2, T] \quad (12)$$

154 and construct a data-driven threshold defined as

$$\mathcal{T}_{\text{thr}} := \text{mean}(\Delta \hat{\boldsymbol{\zeta}}) + \mathcal{Z}_q \times \text{std}(\Delta \hat{\boldsymbol{\zeta}}) \quad (13)$$

155 where \mathcal{Z}_q is the $q\%$ quantile of standard Normal distribution. Finally, we declare a change point C_k when
 156 $\Delta \hat{\boldsymbol{\zeta}}^{C_k} > \mathcal{T}_{\text{thr}}$.

157 The data-driven threshold in (13) is intuitive, as the standardized differences $\Delta\hat{\zeta}$ between two consecutive
 158 change points are close to zeros, while the differences that are at the change points are substantially greater
 159 than zeros. When traced in a plot over time t , the $\Delta\hat{\zeta}$ can exhibit the magnitude of structural changes, and
 160 the threshold that deviates from the mean provides a reasonable cut-off value for the standardized differences,
 161 as demonstrated in Figures 5 and 6. In summary, the localizing method derived from the prior distribution
 162 has a statistical justification, while the localizing method with the data-driven threshold is more robust for
 163 different types of network data in practice.

164 4.2 Model Selection

165 The optimization problem in (2) involves a tuning parameter that can yield different sets of detected change
 166 points when it is varied. In this work, we use Cross-Validation to select λ . In particular, we split the original
 167 time series of graphs into training and testing sets: the training set consists of graphs at odd indexed time
 168 points and the testing set consists of graphs at even indexed time points. Fixed on a specific λ value, we
 169 learn the model parameters with the training set, and we evaluate the learned model with the testing set.

For a list of λ values, we choose the λ giving the maximal log-likelihood on the testing set. Note that the
 log-likelihood is approximated by Monte Carlo samples $\{\mathbf{z}_u^t\}_{u=1}^s$ drawn from the prior distribution $P(\mathbf{z}^t)$ as

$$\sum_{t=1}^T \log P(\mathbf{y}^t) \approx \sum_{t=1}^T \log \left[\frac{1}{s} \sum_{u=1}^s \left[\prod_{(i,j) \in \mathcal{Y}} P_{\phi}(\mathbf{y}_{ij}^t | \mathbf{z}_u^t) \right] \right].$$

170 Further computational details are discussed in Appendix 7.4. Anchored on the selected λ value, we learn
 171 the model parameters again with the full data, resulting the final set of detected change points.

172 5 Simulated and Real Data Experiments

In this section, we implement the proposed method on simulated and real data. To evaluate the performance
 for simulated data, we use three standard metrics in the literature that focus on the number of change points,
 the time gap between the true and detected change points, and the coverage between the segmented time
 intervals. The first metric is the absolute error $|\hat{K} - K|$, where \hat{K} and K are the respective numbers of
 the detected and true change points. The second metric described in Madrid Padilla et al. (2021) is the
 one-sided Hausdorff distance, which is defined as

$$d(\hat{\mathcal{C}}|\mathcal{C}) = \max_{c \in \mathcal{C}} \min_{\hat{c} \in \hat{\mathcal{C}}} |\hat{c} - c|$$

where $\hat{\mathcal{C}}$ and \mathcal{C} are the respective sets of detected and true change points. Also, we report the reversed
 one-sided Hausdorff distance $d(\mathcal{C}|\hat{\mathcal{C}})$. By convention, when $\hat{\mathcal{C}} = \emptyset$, we let $d(\hat{\mathcal{C}}|\mathcal{C}) = \infty$ and $d(\mathcal{C}|\hat{\mathcal{C}}) = -\infty$.
 The last metric described in van den Burg & Williams (2020) is the coverage of a partition \mathcal{G} by another
 partition \mathcal{G}' , which is defined as

$$C(\mathcal{G}, \mathcal{G}') = \frac{1}{T} \sum_{\mathcal{A} \in \mathcal{G}} |\mathcal{A}| \cdot \max_{\mathcal{A}' \in \mathcal{G}'} \frac{|\mathcal{A} \cap \mathcal{A}'|}{|\mathcal{A} \cup \mathcal{A}'|}$$

173 with $\mathcal{A}, \mathcal{A}' \subseteq [1, T]$. The \mathcal{G} and \mathcal{G}' are collections of intervals between consecutive change points for the
 174 respective true and detected change points.

175 5.1 Simulation Study

176 We simulate dynamic graphs from three scenarios to compare the performance of the proposed and competing
 177 methods: Separable Temporal Exponential Random Graph Model, Stochastic Block Model, and Recurrent
 178 Neural Network. For each scenario with different numbers of nodes $n \in \{50, 100\}$, we simulate 10 Monte
 179 Carlo trials of directed dynamic graphs with time span $T = 100$. The true change points are located at
 180 $t = 26, 51, 76$, so the number of change points $K = 3$. Moreover, the $K + 1 = 4$ intervals in the partition \mathcal{G}

181 are $\mathcal{A}_1 = [1, \dots, 25]$, $\mathcal{A}_2 = [26, \dots, 50]$, $\mathcal{A}_3 = [51, \dots, 75]$, and $\mathcal{A}_4 = [76, \dots, 100]$. In each specification, we
 182 report the means and standard deviations over 10 Monte Carlo trials for the evaluation metrics. CPDlatent_N
 183 denotes our proposed approach with the data-driven threshold in (13), using 90% quantile from standard
 184 Normal distribution. We let the latent dimensions $d = 10$ and $k = 5$ for the graph decoder. CPDlatent_G
 185 denotes our proposed approach with the localizing method in (11), using $\alpha = 0.01$ from Gamma distribution.
 186 We let the latent dimensions $k = 10$ and $d = n/10$ for the graph decoder. The number of samples drawn
 187 from the Gamma distribution is $m = 500$ when $d = 10$ and $m = 1000$ when $d = 5$.

188 Three competitors, gSeg (Chen & Zhang, 2015), kerSeg (Song & Chen, 2022b), and CPDstergm (Kei et al.,
 189 2023), are provided for comparison. The gSeg utilizes graph-based scan statistics and kerSeg employs a kernel-
 190 based framework to test the partition before and after a change point. The CPDstergm fits a STERGM with
 191 user-specified network statistics to detect change points based on the model parameters. For CPDstergm,
 192 we first use two network statistics, edge count and mutuality, in both formation and dissolution models to let
 193 $p = 4$. We then add one more network statistic, number of triangles, to let $p = 6$ as another specification. For
 194 gSeg, we use the minimum spanning tree to construct the similarity graph, with the approximated p-value
 195 of the original edge-count scan statistic, and we set $\alpha = 0.05$. For kerSeg, we use the approximated p-value
 196 of the fGKCP₁, and we set $\alpha = 0.001$. Moreover, we use networks (nets.) and network statistics (stats.) as
 197 two types of input data to gSeg and kerSeg. Throughout, we choose these settings because they produce
 198 good performance on average for the competitors. Changing the settings can enhance their performance on
 199 some specifications, while severely jeopardizing their performance on other specifications.

200 Scenario 1: Separable Temporal Exponential Random Graph Model

In this scenario, we apply time-homogeneous Separable Temporal Exponential Random Graph Model (STERGM) between change points to generate sequences of dynamic networks (Krivitsky & Handcock, 2014). We use three network statistics, edge count, mutuality, and number of triangles, in both formation (F) and dissolution (D) models. The $p = 6$ parameters for each time point t are

$$\theta_F^t, \theta_D^t = \begin{cases} -2, 2, -2, -1, 2, 1, & t \in \mathcal{A}_1 \cup \mathcal{A}_3 \setminus 1, \\ -1.5, 1, -1, 2, 1, 1.5, & t \in \mathcal{A}_2 \cup \mathcal{A}_4. \end{cases}$$

201 Figure 2 exhibits examples of generated networks. Visually, STERGM produces adjacency matrices that are
 202 sparse, which is often the case in real world social networks.

203 Table 1 displays the means and standard deviations of the evaluation metrics for comparison. Since the
 204 networks are directly sampled from STERGM, the CPDstergm method with correctly specified network
 205 statistics ($p = 6$) achieves the best result, in terms of greater converge of the intervals. However, when
 206 the network statistics are mis-specified ($p = 4$), the performance of CPDstergm is worsened, with greater
 207 gaps between the true and detected change points. Also, using either networks (nets.) or network statistics
 208 (stats.) cannot improve the performance of gSeg and kerSeg methods: the binary search approach tend to
 209 detect excessive number of change points. Our CPDlatent method, without the need of specifying network
 210 statistics, can achieve relatively good performance on average.

211 Scenario 2: Stochastic Block Model

In this scenario, we use Stochastic Block Model (SBM) to generate sequences of dynamic networks, and we impose a time-dependent mechanism in the generation process as in Madrid Padilla et al. (2022). Two probability matrices $\mathbf{P}, \mathbf{Q} \in [0, 1]^{n \times n}$ are constructed and they are defined as

$$\mathbf{P}_{ij} = \begin{cases} 0.5, & i, j \in \mathcal{B}_l, l \in [3], \\ 0.3, & \text{otherwise,} \end{cases} \quad \text{and} \quad \mathbf{Q}_{ij} = \begin{cases} 0.45, & i, j \in \mathcal{B}_l, l \in [3], \\ 0.2, & \text{otherwise,} \end{cases}$$

where $\mathcal{B}_1, \mathcal{B}_2, \mathcal{B}_3$ are evenly sized clusters that form a partition of $\{1, \dots, n\}$. Then a sequence of matrices $\mathbf{E}^t \in [0, 1]^{n \times n}$ are arranged for $t = 1, \dots, T$ such that

$$\mathbf{E}_{ij}^t = \begin{cases} \mathbf{P}_{ij}, & t \in \mathcal{A}_1 \cup \mathcal{A}_3, \\ \mathbf{Q}_{ij}, & t \in \mathcal{A}_2 \cup \mathcal{A}_4. \end{cases}$$

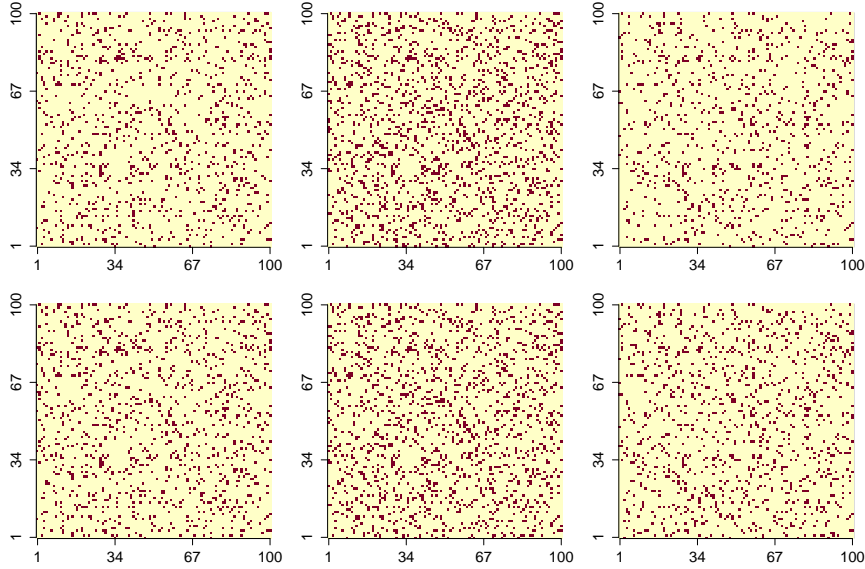


Figure 2: Examples of networks generated from STERGM with $n = 100$. In the first row, from left to right, each plot corresponds to the network at $t = 25, 50, 75$ respectively. In the second row, from left to right, each plot corresponds to the network at $t = 26, 51, 76$ respectively (the change points).

Table 1: Means (standard deviations) of evaluation metrics for dynamic graphs simulated from STERGM. The best coverage metric is bolded.

n	Method	$ \hat{K} - K \downarrow$	$d(\hat{\mathcal{C}} \mathcal{C}) \downarrow$	$d(\mathcal{C} \hat{\mathcal{C}}) \downarrow$	$C(\mathcal{G}, \mathcal{G}') \uparrow$
50	CPDlatent $_N$	0.1 (0.3)	4.3 (5.7)	2.6 (1.3)	90.87%
	CPDlatent $_G$	0.4 (0.6)	4.2 (6.9)	3.4 (3.4)	90.97%
	CPDstergm $_{p=4}$	1.5 (0.8)	11.7 (7.5)	10.5 (2.3)	67.68%
	CPDstergm $_{p=6}$	0.2 (0.4)	1.6 (1.2)	3 (3.5)	91.54%
	gSeg (nets.)	12.3 (0.5)	0 (0)	19 (0)	27.90%
	gSeg (stats.)	15.8 (0.7)	1.5 (0.5)	20.1 (0.3)	24.55%
	kerSeg (nets.)	9.7 (0.9)	1.4 (0.9)	17.9 (1.2)	37.62%
	kerSeg (stats.)	9.4 (0.7)	3.9 (1.3)	18 (1.8)	35.86%
100	CPDlatent $_N$	0 (0)	3.9 (1.3)	3.9 (1.3)	91.33%
	CPDlatent $_G$	0.7 (1.3)	3.1 (1.3)	6.0 (4.0)	88.55%
	CPDstergm $_{p=4}$	0.7 (0.6)	21.9 (10.3)	7.6 (4.3)	67.21%
	CPDstergm $_{p=6}$	0 (0)	1.1 (0.3)	1.1 (0.3)	94.01%
	gSeg (nets.)	12 (0)	0 (0)	19 (0)	28.00%
	gSeg (stats.)	14.5 (2.3)	3.3 (3.6)	20.2 (0.4)	26.13%
	kerSeg (nets.)	9.3 (0.8)	1 (0)	17.7 (0.6)	37.62%
	kerSeg (stats.)	8.5 (0.8)	4.5 (1.4)	17.3 (1.7)	36.92%

Lastly, the networks are generated with $\rho = 0.5$ as a time-dependent mechanism. For $t = 1, \dots, T - 1$, we let $\mathbf{y}_{ij}^1 \sim \text{Bernoulli}(\mathbf{E}_{ij}^1)$ and

$$\mathbf{y}_{ij}^{t+1} \sim \begin{cases} \text{Bernoulli}(\rho(1 - \mathbf{E}_{ij}^{t+1}) + \mathbf{E}_{ij}^{t+1}), & \mathbf{y}_{ij}^t = 1, \\ \text{Bernoulli}((1 - \rho)\mathbf{E}_{ij}^{t+1}), & \mathbf{y}_{ij}^t = 0. \end{cases}$$

212 With $\rho > 0$, the probability to form an edge for i, j becomes greater at time $t + 1$ when there exists an edge
 213 at time t , and the probability becomes smaller when there does not exist an edge at time t . Figure 3 exhibits
 214 examples of generated networks. Visually, SBM produces adjacency matrices with block structures, where
 215 mutuality serves as an important pattern for the homophily within groups.

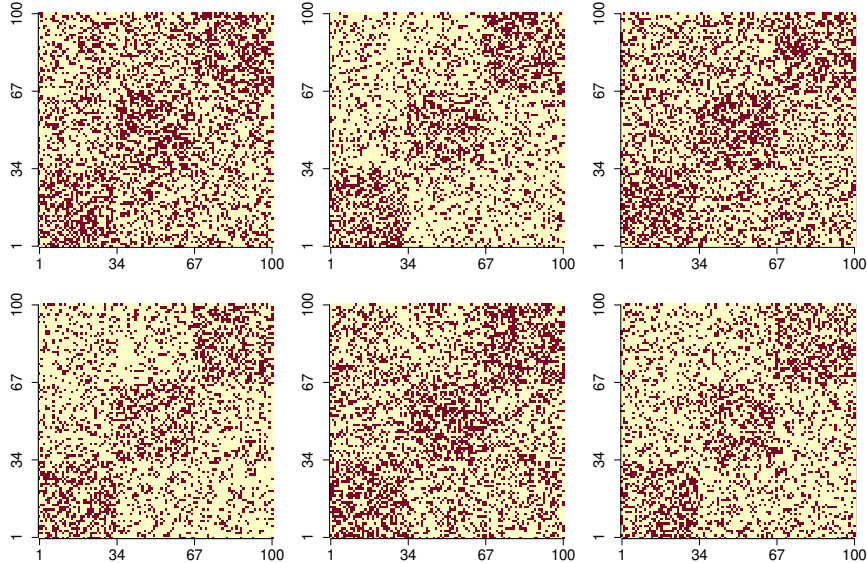


Figure 3: Examples of networks generated from SBM with $n = 100$. In the first row, from left to right, each plot corresponds to the network at $t = 25, 50, 75$ respectively. In the second row, from left to right, each plot corresponds to the network at $t = 26, 51, 76$ respectively (the change points).

216 Table 2 displays the means and standard deviations of the evaluation metrics for comparison. As expected,
 217 both CPDstergm methods with $p = 4$ and $p = 6$ that utilize the mutuality as a network statistic for the
 218 detection can achieve good results, in terms of greater converge of the intervals. Furthermore, using network
 219 statistics (stats.) for both gSeg and kerSeg methods can improve their performance, comparing to using
 220 networks (nets.) as input data. Lastly, our CPDlatent method, which infers the features in latent space that
 221 induce the structural changes, achieves the best result for networks with block structures.

222 Scenario 3: Recurrent Neural Networks

In this scenario, we use Recurrent Neural Networks (RNN) to generate sequences of dynamic networks. Specifically, we sample latent variables from pre-defined priors, and we initialize the RNN with uniform weights. The graphs are then generated by the matrix multiplication defined in Section 2.2, using the output of RNN. The parameters for the pre-defined priors are

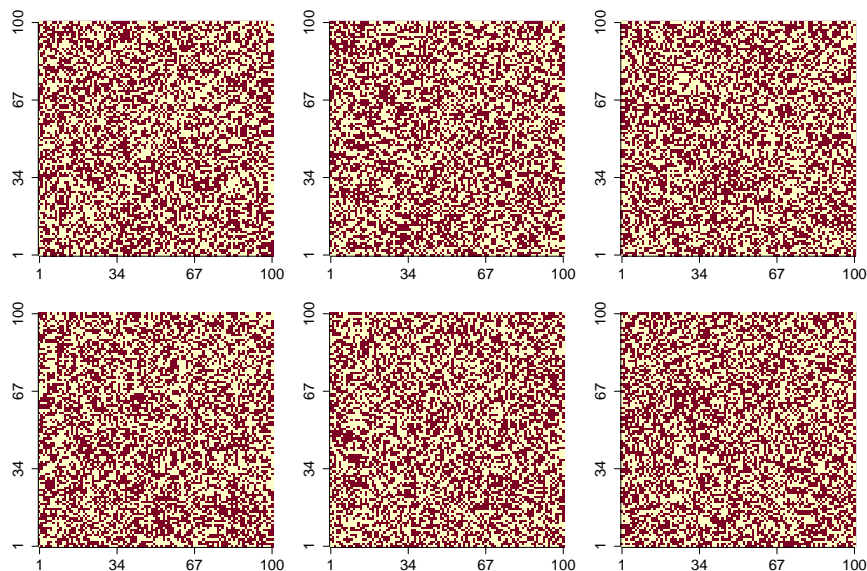
$$z^t \sim \begin{cases} \mathcal{N}(-\mathbf{1}, 0.1\mathbf{I}_d), & t \in \mathcal{A}_1 \cup \mathcal{A}_3, \\ \mathcal{N}(\mathbf{5}, 0.1\mathbf{I}_d), & t \in \mathcal{A}_2 \cup \mathcal{A}_4. \end{cases}$$

223 Similar to the previous two scenarios, the simulation using RNN also imposes a time-dependent mechanism
 224 across dynamic networks. Figure 4 exhibits examples of generated networks. Visually, RNN produces
 225 adjacency matrices that are dense, and no discernible pattern can be noticed.

226 Table 3 displays the means and standard deviations of the evaluation metrics for comparison. Because no
 227 structural pattern or suitable network statistics can be determined a priori, neither CPDstergm method with
 228 $p = 4$ nor with $p = 6$ can detect the change points accurately. Likewise, both gSeg and kerSeg methods that
 229 utilize the mis-specified network statistics (stats.) cannot produce satisfactory performance. Notably, the
 230 kerSeg method that exploits the features in high dimension with networks (nets.) instead of user-specified

Table 2: Means (stds.) of evaluation metrics for dynamic networks simulated from SBM. The best coverage metric is bolded.

n	Method	$ \hat{K} - K \downarrow$	$d(\hat{\mathcal{C}} \mathcal{C}) \downarrow$	$d(\mathcal{C} \hat{\mathcal{C}}) \downarrow$	$C(\mathcal{G}, \mathcal{G}') \uparrow$
50	CPDlatent $_N$	0 (0)	0.1 (0.3)	0.1 (0.3)	99.80%
	CPDlatent $_G$	0.3 (0.6)	0.1 (0.3)	3.1 (6.2)	96.70%
	CPDstergm $_{p=4}$	0.1 (0.3)	1 (0)	2.4 (4.2)	97.04%
	CPDstergm $_{p=6}$	0.3 (0.5)	1 (0)	4.6 (5.6)	94.74%
	gSeg (nets.)	12.9 (1.8)	0 (0)	19.4 (0.8)	27.20%
	gSeg (stats.)	2.2 (0.7)	Inf (na)	-Inf (na)	49.21%
	kerSeg (nets.)	6.4 (1.4)	0 (0)	16.6 (2.0)	45.50%
	kerSeg (stats.)	0.9 (1.2)	0 (0)	5.6 (6.8)	93.50%
100	CPDlatent $_N$	0.1 (0.3)	0.1 (0.3)	1.3 (3.6)	98.60%
	CPDlatent $_G$	0.5 (0.7)	0.2 (0.4)	5.1 (6.1)	94.81%
	CPDstergm $_{p=4}$	0 (0)	1 (0)	1 (0)	98.04%
	CPDstergm $_{p=6}$	0 (0)	1 (0)	1 (0)	98.04%
	gSeg (nets.)	12.3 (0.9)	0 (0)	19 (0)	27.80%
	gSeg (stats.)	2 (0.4)	Inf (na)	-Inf (na)	55.75%
	kerSeg (nets.)	6 (0.8)	0 (0)	15.2 (2.0)	47.00%
	kerSeg (stats.)	0.9 (0.7)	0 (0)	9.6 (7.6)	93.40%

Figure 4: Examples of networks generated from RNN with $n = 100$. In the first row, from left to right, each plot corresponds to the network at $t = 25, 50, 75$ respectively. In the second row, from left to right, each plot corresponds to the network at $t = 26, 51, 76$ respectively (the change points).

231 network statistics (stats.) can deliver a good result. Lastly, our CPDlatent method that first infers the
232 graph level representations from the complex network structures and then utilize them to detect the change
233 points yields the best result.

Table 3: Means (stds.) of evaluation metrics for dynamic networks simulated from RNN. The best coverage metric is bolded.

n	Method	$ \hat{K} - K \downarrow$	$d(\hat{\mathcal{C}} \mathcal{C}) \downarrow$	$d(\mathcal{C} \hat{\mathcal{C}}) \downarrow$	$C(\mathcal{G}, \mathcal{G}') \uparrow$
50	CPDlatent _{N}	0 (0)	1.8 (0.7)	1.8 (0.7)	94.77%
	CPDlatent _{G}	0.3 (0.6)	1.7 (0.6)	3.2 (3.0)	93.04%
	CPDstergm _{$p=4$}	2.0 (1.7)	6.0 (7.7)	15.2 (4.9)	72.10%
	CPDstergm _{$p=6$}	1.0 (0.4)	18.5 (9.4)	14.3 (2.9)	60.25%
	gSeg (nets.)	2.3 (0.6)	Inf (na)	-Inf (na)	29.42%
	gSeg (stats.)	2.9 (0.3)	Inf (na)	-Inf (na)	2.47%
	kerSeg (nets.)	1.5 (0.9)	1.4 (0.7)	5.3 (3.3)	89.25%
	kerSeg (stats.)	2.8 (0.4)	Inf (na)	-Inf (na)	9.89%
100	CPDlatent _{N}	0 (0)	2.5 (0.7)	2.5 (0.7)	91.96%
	CPDlatent _{G}	0.2 (0.6)	2.1 (0.7)	2.8 (1.8)	92.34%
	CPDstergm _{$p=4$}	2.0 (1.4)	10.6 (8.0)	14.1 (3.1)	60.37%
	CPDstergm _{$p=6$}	1.2 (1.3)	20.6 (12.6)	15.2 (5.9)	53.21%
	gSeg (nets.)	3 (0)	Inf (na)	-Inf (na)	0%
	gSeg (stats.)	2.9 (0.3)	Inf (na)	-Inf (na)	4.27%
	kerSeg (nets.)	1.4 (0.7)	1.9 (0.7)	5.4 (1.9)	88.95%
	kerSeg (stats.)	3 (0)	Inf (na)	-Inf (na)	0%

234 5.2 MIT Cellphone Data

235 The Massachusetts Institute of Technology (MIT) cellphone data (Eagle & Pentland, 2006) depicts human
 236 interactions via phone call activities among $n = 96$ participants spanning $T = 232$ days. An edge $\mathbf{y}_{ij}^t = 1$ in
 237 the constructed networks indicates that participant i and participant j had made phone calls on day t , and
 238 $\mathbf{y}_{ij}^t = 0$ otherwise. The data ranges from 2004-09-15 to 2005-05-04, covering the winter break in the MIT
 239 academic calendar.

240 We detect the change points with our proposed method using the data-driven threshold from standard
 241 Normal distribution, and we use network statistics as input data to the competitor methods. Specifically,
 242 we use the number of edges, isolates, and triangles to capture the frequency of connections, the sparsity
 243 of social interaction, and the transitive association among friends, respectively. Figure 5 displays $\Delta\hat{\zeta}$ of
 244 Equation (12), and the detected change points from our method and competitor methods. Furthermore,
 245 Table 4 provides a list of potential events, aligning with the detected change points from our method.

246 Without specifying the structural changes to search for, our method can punctually detect the beginning
 247 of the winter break, which is the major event that alters the interaction among participants. Similar to
 248 the competitors, we have detected a spike on 2004-10-23, corresponding to the annual sponsor meeting that
 249 occurred on 2004-10-21. More than two-thirds of the participants have attended the meeting, focusing on
 250 achieving project goals throughout the week (Eagle & Pentland, 2006). Moreover, we have detected other
 251 change points related to national holidays and spring break.

252 5.3 Enron Email Data

253 The Enron email data, analyzed by Priebe et al. (2005), Park et al. (2012), and Peel & Clauset (2015), por-
 254 trays communication among employees before the collapse of a giant energy company. The dynamic network
 255 data consists of $T = 100$ weekly networks, ranging from 2000-06-05 to 2002-05-06 for $n = 100$ employees. We
 256 detect the change points with our proposed method using the data-driven threshold from standard Normal
 257 distribution, and we use the same network statistics described in Section 5.2 to the competitor methods.

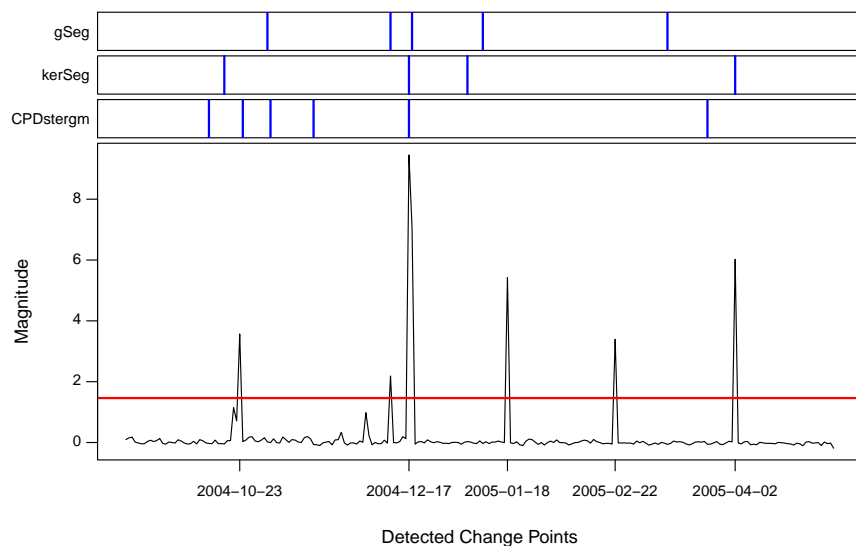


Figure 5: Detected change points from the proposed and competitor (blue) methods on the MIT Cellphone Data. The threshold (red horizontal line) is calculated by (13) with $\mathcal{Z}_{0.9}$.

Table 4: Potential nearby events aligned with the detected change points from our proposed method on the MIT cellphone data.

Detected change points	Potential nearby events
2004-10-23	2004-10-21 Sponsor meeting
2004-12-17	2004-12-18 to 2005-01-02 Winter break
2005-01-18	2005-01-17 Martin Luther King Day
2005-02-22	2005-02-21 Presidents Day
2005-04-02	2005-03-21 to 2005-03-25 Spring break

258 Figure 6 displays $\Delta\hat{\zeta}$ of Equation (12), and the detected change points from our method and competitor
 259 methods. Furthermore, Table 5 provides a list of potential events, aligning with the detected change points
 260 from our method.

261 In 2001, Enron underwent a multitude of major and overlapping incidents, making it difficult to associate
 262 the detected change points with specific real world events. Yet, as our proposed method detects the results
 263 over the two-year time frame, four crucial change points are detected for interpretation. Throughout 2000,
 264 Enron orchestrated rolling blackouts, causing staggering surges in electricity prices that peaked at twenty
 265 times the standard rate. The situation worsened when the Federal Energy Regulatory Commission (FERC)
 266 exonerated Enron of wrongdoing by the end of 2000. During a public appearance in June 2001, the CEO is
 267 physically confronted by an activist in protest against Enron’s role in the energy crisis. Amid the turmoil,
 268 an employee meeting took place in September 2001, where the CEO reassured employees that Enron’s stock
 269 was a good buy and the company’s accounting methods were legal and appropriate. Following the employee
 270 meeting, the stock saw a brief surge before continuing its sharp decline. Three months later, pressured by
 271 Wall Street analysts and the revelation of the scandals, Enron filed for bankruptcy and the largest energy
 272 company in the U.S. fell apart.

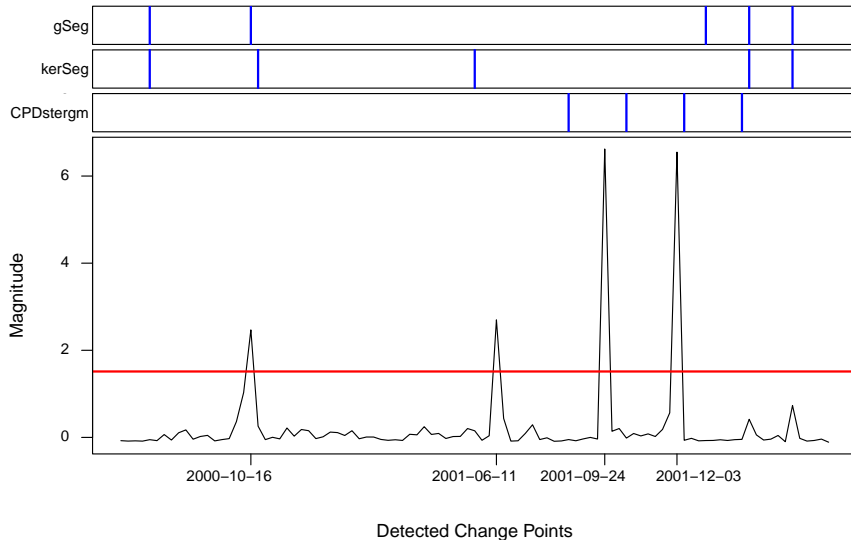


Figure 6: Detected change points from the proposed and competitor (blue) methods on the Enron email data. The threshold (red horizontal line) is calculated by (13) with $\mathcal{Z}_{0.9}$.

Table 5: Potential nearby events aligned with the detected change points from our proposed method on the Enron email data.

Detected change points	Potential nearby events
2000-10-16	2000-11-01 FERC exonerated Enron
2001-06-11	2001-06-21 CEO publicly confronted
2001-09-24	2001-09-26 Employee meeting
2001-12-03	2001-12-02 Enron filed for bankruptcy

6 Discussion

273

274 This paper proposes a generative model to detect change points in dynamic graphs. Intrinsically, dynamic
 275 networks are complex due to dyadic and temporal dependencies. Learning low dimensional graph represen-
 276 tations can extract useful features to facilitate change point detection in dynamic graphs. We impose prior
 277 distributions to the graph representations, and the priors for the latent space are learned from the data
 278 as empirical Bayes. The optimization problem with Group Fused Lasso penalty is solved via ADMM, and
 279 generative model is demonstrated to be useful for change point detection.

280 Several extensions to our proposed framework are possible for future development. Besides binary networks,
 281 relations by nature have degree of strength, which are denoted by generic values. Also, nodal and dyadic
 282 attributes are important components in network data. Hence, models that can generate weighted edges, as
 283 well as nodal and dyadic attributes, can capture more information about the network dynamics (Fellows &
 284 Hancock, 2012; Krivitsky, 2012; Simonovsky & Komodakis, 2018). Furthermore, the number of nodes and
 285 their attributes are subjected to change over time. Extending the framework to allow the network size to
 286 change and to detect vertex level anomalies can provide granular insights in addition to graph level changes
 287 (Simonovsky & Komodakis, 2018; Shen et al., 2023). Similarly, improving the scalability and computational
 288 efficiency for representation learning is also crucial (Killick et al., 2012; Gallagher et al., 2021), especially
 289 for handling large and weighted graphs. While our framework demonstrates the ability in change point
 290 detection, the development of more sophisticated architectures can enhance the model’s capacity on other
 291 meaningful tasks (Hancock et al., 2007; Kolar et al., 2010; Yu et al., 2021; Madrid Padilla et al., 2023).

References

- 292 Carlos M Alaíz, Alvaro Barbero, and José R Dorronsoro. Group fused lasso. In *Artificial Neural Networks*
293 *and Machine Learning–ICANN 2013: 23rd International Conference on Artificial Neural Networks Sofia,*
294 *Bulgaria, September 10–13, 2013. Proceedings 23*, pp. 66–73. Springer, 2013.
- 296 Kevin Bleakley and Jean-Philippe Vert. The group fused lasso for multiple change-point detection. *arXiv*
297 *preprint arXiv:1106.4199*, 2011.
- 298 Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization
299 and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in*
300 *Machine Learning*, 3(1):1–122, 2011.
- 301 Carter T Butts. A relational event framework for social action. *Sociological Methodology*, 38(1):155–200,
302 2008.
- 303 Carter T Butts, Alessandro Lomi, Tom AB Snijders, and Christoph Stadtfeld. Relational event models in
304 network science. *Network Science*, 11(2):175–183, 2023.
- 305 Guodong Chen, Jesús Arroyo, Avanti Athreya, Joshua Cape, Joshua T Vogelstein, Youngser Park, Chris
306 White, Jonathan Larson, Weiwei Yang, and Carey E Priebe. Multiple network embedding for anomaly
307 detection in time series of graphs. *arXiv preprint arXiv:2008.10055*, 2020.
- 308 Hao Chen and Nancy Zhang. Graph-based change-point detection. *The Annals of Statistics*, 43(1):139–176,
309 2015.
- 310 Lynna Chu and Hao Chen. Asymptotic distribution-free change-point detection for multivariate and non-
311 euclidean data. *The Annals of Statistics*, 47(1):382–414, 2019.
- 312 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional
313 transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 314 Vijay Prakash Dwivedi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Graph neural
315 networks with learnable structural and positional representations. *arXiv preprint arXiv:2110.07875*, 2021.
- 316 Nathan Eagle and Alex (Sandy) Pentland. Reality mining: sensing complex social systems. *Personal and*
317 *Ubiquitous Computing*, 10(4):255–268, 2006.
- 318 Ian Fellows and Mark S. Handcock. Exponential-family random network models, 2012.
- 319 Ian Gallagher, Andrew Jones, and Patrick Rubin-Delanchy. Spectral embedding for dynamic networks with
320 stability guarantees. *Advances in Neural Information Processing Systems*, 34:10158–10170, 2021.
- 321 Damien Garreau and Sylvain Arlot. Consistent change-point detection with kernels. *Electronic Journal of*
322 *Statistics*, 12(2):4440 – 4486, 2018.
- 323 Yongshun Gong, Xue Dong, Jian Zhang, and Meng Chen. Latent evolution model for change point detection
324 in time-varying networks. *Information Sciences*, 646:119376, 2023.
- 325 Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances*
326 *in neural information processing systems*, 30, 2017.
- 327 Mingqi Han, Eric A Bushong, Mayuko Segawa, Alexandre Tiard, Alex Wong, Morgan R Brady, Milica
328 Momcilovic, Dane M Wolf, Ralph Zhang, Anton Petcherski, et al. Spatial mapping of mitochondrial
329 networks and bioenergetics in lung cancer. *Nature*, 615(7953):712–719, 2023.
- 330 Mark S Handcock, Adrian E Raftery, and Jeremy M Tantrum. Model-based clustering for social networks.
331 *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 170(2):301–354, 2007.
- 332 Steve Hanneke, Wenjie Fu, and Eric P Xing. Discrete temporal models of social networks. *Electronic Journal*
333 *of Statistics*, 4:585–605, 2010.

- 334 Xiaoxin He, Bryan Hooi, Thomas Laurent, Adam Perold, Yann LeCun, and Xavier Bresson. A generalization
335 of vit/mlp-mixer to graphs. In *International Conference on Machine Learning*, pp. 12724–12745. PMLR,
336 2023.
- 337 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural
338 information processing systems*, 33:6840–6851, 2020.
- 339 Shenyang Huang, Yasmeeen Hitti, Guillaume Rabusseau, and Reihaneh Rabbany. Laplacian change point
340 detection for dynamic graphs. In *Proceedings of the 26th ACM SIGKDD International Conference on
341 Knowledge Discovery & Data Mining*, pp. 349–358, 2020.
- 342 Yik Lun Kei, Hangjian Li, Yanzhen Chen, and Oscar Hernan Madrid Padilla. Change point detection on a
343 separable model for dynamic networks. *arXiv preprint arXiv:2303.17642*, 2023.
- 344 Rebecca Killick, Paul Fearnhead, and Idris A Eckley. Optimal detection of changepoints with a linear
345 computational cost. *Journal of the American Statistical Association*, 107(500):1590–1598, 2012.
- 346 Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- 347 Thomas N Kipf and Max Welling. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*, 2016.
- 348 Mladen Kolar, Le Song, Amr Ahmed, and Eric P Xing. Estimating time-varying networks. *The Annals of
349 Applied Statistics*, pp. 94–123, 2010.
- 350 Pavel N Krivitsky. Exponential-family random graph models for valued networks. *Electronic journal of
351 statistics*, 6:1100, 2012.
- 352 Pavel N Krivitsky and Mark S Handcock. A separable model for dynamic networks. *Journal of the Royal
353 Statistical Society. Series B, Statistical Methodology*, 76(1):29, 2014.
- 354 Federico Larroca, Paola Bermolen, Marcelo Fiori, and Gonzalo Mateos. Change point detection in weighted
355 and directed random dot product graphs. In *2021 29th European Signal Processing Conference (EU-
356 SIPCO)*, pp. 1810–1814. IEEE, 2021.
- 357 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves
358 Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language
359 generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- 360 Carlos Misael Madrid Padilla, Haotian Xu, Daren Wang, Oscar Hernan Madrid Padilla, and Yi Yu. Change
361 point detection and inference in multivariable nonparametric models under mixing conditions. *arXiv
362 preprint arXiv:2301.11491*, 2023.
- 363 Oscar Hernan Madrid Padilla, Yi Yu, Daren Wang, and Alessandro Rinaldo. Optimal nonparametric multi-
364 variate change point detection and localization. *IEEE Transactions on Information Theory*, 68(3):1922–
365 1944, 2021.
- 366 Oscar Hernan Madrid Padilla, Yi Yu, and Carey E Priebe. Change point localization in dependent dynamic
367 nonparametric random dot product graphs. *The Journal of Machine Learning Research*, 23(1):10661–
368 10719, 2022.
- 369 Bernardo Marenco, Paola Bermolen, Marcelo Fiori, Federico Larroca, and Gonzalo Mateos. Online change
370 point detection for weighted and directed random dot product graphs. *IEEE Transactions on Signal and
371 Information Processing over Networks*, 8:144–159, 2022.
- 372 Erik Nijkamp, Mitch Hill, Tian Han, Song-Chun Zhu, and Ying Nian Wu. On the anatomy of mcmc-based
373 maximum likelihood learning of energy-based models. In *Proceedings of the AAAI Conference on Artificial
374 Intelligence*, volume 34, pp. 5272–5280, 2020.
- 375 Shirui Pan, Ruiqi Hu, Guodong Long, Jing Jiang, Lina Yao, and Chengqi Zhang. Adversarially regularized
376 graph autoencoder for graph embedding. *arXiv preprint arXiv:1802.04407*, 2018.

- 377 Bo Pang, Tian Han, Erik Nijkamp, Song-Chun Zhu, and Ying Nian Wu. Learning latent space energy-based
378 prior model. *Advances in Neural Information Processing Systems*, 33:21994–22008, 2020.
- 379 Jong Hee Park and Yunkyu Sohn. Detecting Structural Changes in Longitudinal Network Data. *Bayesian*
380 *Analysis*, 15(1):133 – 157, 2020.
- 381 Youngser Park, Carey E Priebe, and Abdou Youssef. Anomaly detection in time series of graphs using fusion
382 of graph invariants. *IEEE journal of selected topics in signal processing*, 7(1):67–75, 2012.
- 383 Leto Peel and Aaron Clauset. Detecting change points in the large-scale structure of evolving networks. In
384 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, 2015.
- 385 Carey E Priebe, John M Conroy, David J Marchette, and Youngser Park. Scan statistics on enron graphs.
386 *Computational & Mathematical Organization Theory*, 11:229–247, 2005.
- 387 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution
388 image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer*
389 *vision and pattern recognition*, pp. 10684–10695, 2022.
- 390 S Golshid Sharifnia and Abbas Saghaei. A statistical approach for social network change detection: an ergm
391 based framework. *Communications in Statistics-Theory and Methods*, 51(7):2259–2280, 2022.
- 392 Cen Cheng Shen, Jonathan Larson, Ha Trinh, Xihan Qin, Youngser Park, and Carey E Priebe. Discovering
393 communication pattern shifts in large-scale labeled networks using encoder embedding and vertex
394 dynamics. *IEEE Transactions on Network Science and Engineering*, 2023.
- 395 Martin Simonovsky and Nikos Komodakis. Graphvae: Towards generation of small graphs using variational
396 autoencoders. In *Artificial Neural Networks and Machine Learning–ICANN 2018: 27th International*
397 *Conference on Artificial Neural Networks, Rhodes, Greece, October 4-7, 2018, Proceedings, Part I 27*, pp.
398 412–422. Springer, 2018.
- 399 Tom AB Snijders. The statistical evaluation of social network dynamics. *Sociological methodology*, 31(1):
400 361–395, 2001.
- 401 Tom AB Snijders, Gerhard G Van de Bunt, and Christian EG Steglich. Introduction to stochastic actor-based
402 models for network dynamics. *Social networks*, 32(1):44–60, 2010.
- 403 Hoseung Song and Hao Chen. Asymptotic distribution-free changepoint detection for data with repeated
404 observations. *Biometrika*, 109(3):783–798, 2022a.
- 405 Hoseung Song and Hao Chen. New kernel-based change-point detection. *arXiv preprint arXiv:2206.01853*,
406 2022b.
- 407 Deborah Sulem, Henry Kenlay, Mihai Cucuringu, and Xiaowen Dong. Graph similarity learning for change-
408 point detection in dynamic networks. *Machine Learning*, pp. 1–44, 2023.
- 409 Gerrit JJ van den Burg and Christopher KI Williams. An evaluation of change point detection algorithms.
410 *arXiv preprint arXiv:2003.06222*, 2020.
- 411 Jean-Philippe Vert and Kevin Bleakley. Fast detection of multiple change-points shared by many signals
412 using group lars. *Advances in Neural Information Processing Systems*, 23, 2010.
- 413 Yu Wang, Wotao Yin, and Jinshan Zeng. Global convergence of admm in nonconvex nonsmooth optimization.
414 *Journal of Scientific Computing*, 78:29–63, 2019.
- 415 Jianwen Xie, Song-Chun Zhu, and Ying Nian Wu. Synthesizing dynamic patterns by spatial-temporal
416 generative convnet. In *Proceedings of the ieee conference on computer vision and pattern recognition*, pp.
417 7093–7101, 2017.

- 418 Jianwen Xie, Yang Lu, Ruiqi Gao, Song-Chun Zhu, and Ying Nian Wu. Cooperative training of descriptor
419 and generator networks. *IEEE transactions on pattern analysis and machine intelligence*, 42(1):27–45,
420 2018.
- 421 Yi Yu, Oscar Hernan Madrid Padilla, Daren Wang, and Alessandro Rinaldo. Optimal network online change
422 point localisation. *arXiv preprint arXiv:2101.05477*, 2021.
- 423 Ming Yuan and Yi Lin. Model selection and estimation in regression with grouped variables. *Journal of the*
424 *Royal Statistical Society Series B: Statistical Methodology*, 68(1):49–67, 2006.
- 425 Xinxun Zhang, Pengfei Jiao, Mengzhou Gao, Tianpeng Li, Yiming Wu, Huaming Wu, and Zhidong Zhao.
426 Vggm: Variational graph gaussian mixture model for unsupervised change point detection in dynamic
427 networks. *IEEE Transactions on Information Forensics and Security*, 2024.
- 428 Zifeng Zhao, Li Chen, and Lizhen Lin. Change-point detection in dynamic networks via graphon estimation.
429 *arXiv preprint arXiv:1908.01823*, 2019.

430 7 Appendix

431 7.1 Updating μ and ϕ

In this section, we derive the updates for prior parameter $\mu \in \mathbb{R}^{T \times d}$ and graph decoder parameter ϕ . Denote the objective function in Equation (4) as $\mathcal{L}(\phi, \mu)$ and denote the set of parameters $\{\phi, \mu\}$ as θ . We first calculate the gradient of the log-likelihood $l(\theta)$ in $\mathcal{L}(\phi, \mu)$ with respect to θ :

$$\begin{aligned}
\nabla_{\theta} l(\theta) &= \nabla_{\theta} \sum_{t=1}^T \log P(\mathbf{y}^t) \\
&= \sum_{t=1}^T \frac{1}{P(\mathbf{y}^t)} \nabla_{\theta} P(\mathbf{y}^t) \\
&= \sum_{t=1}^T \frac{1}{P(\mathbf{y}^t)} \nabla_{\theta} \int P(\mathbf{y}^t, \mathbf{z}^t) d\mathbf{z}^t \\
&= \sum_{t=1}^T \frac{1}{P(\mathbf{y}^t)} \int P(\mathbf{y}^t, \mathbf{z}^t) [\nabla_{\theta} \log P(\mathbf{y}^t, \mathbf{z}^t)] d\mathbf{z}^t \\
&= \sum_{t=1}^T \int \frac{P(\mathbf{y}^t, \mathbf{z}^t)}{P(\mathbf{y}^t)} [\nabla_{\theta} \log P(\mathbf{y}^t, \mathbf{z}^t)] d\mathbf{z}^t \\
&= \sum_{t=1}^T \int P(\mathbf{z}^t | \mathbf{y}^t) [\nabla_{\theta} \log P(\mathbf{y}^t, \mathbf{z}^t)] d\mathbf{z}^t \\
&= \sum_{t=1}^T \mathbb{E}_{P(\mathbf{z}^t | \mathbf{y}^t)} \left(\nabla_{\theta} \log [P(\mathbf{y}^t | \mathbf{z}^t) P(\mathbf{z}^t)] \right) \\
&= \sum_{t=1}^T \mathbb{E}_{P(\mathbf{z}^t | \mathbf{y}^t)} \left(\nabla_{\theta} \log P(\mathbf{y}^t | \mathbf{z}^t) \right) + \sum_{t=1}^T \mathbb{E}_{P(\mathbf{z}^t | \mathbf{y}^t)} \left(\nabla_{\theta} \log P(\mathbf{z}^t) \right).
\end{aligned}$$

Note that the expectation in the gradient is now with respect to the posterior distribution $P(\mathbf{z}^t | \mathbf{y}^t) \propto P(\mathbf{y}^t | \mathbf{z}^t) \times P(\mathbf{z}^t)$. Furthermore, the gradient of $\mathcal{L}(\phi, \mu)$ with respect to the prior parameter $\mu^t \in \mathbb{R}^d$ at a specific time point t is

$$\begin{aligned}
\nabla_{\mu^t} \mathcal{L}(\phi, \mu) &= -\mathbb{E}_{P(\mathbf{z}^t | \mathbf{y}^t)} \left(\nabla_{\mu^t} \log P(\mathbf{z}^t) \right) + \kappa(\mu^t - \nu^t + \mathbf{w}^t) \\
&= -\mathbb{E}_{P(\mathbf{z}^t | \mathbf{y}^t)} (\mathbf{z}^t - \mu^t) + \kappa(\mu^t - \nu^t + \mathbf{w}^t).
\end{aligned}$$

Setting the gradient $\nabla_{\mu^t} \mathcal{L}(\phi, \mu)$ to zeros and solve for μ^t , we have

$$\begin{aligned}
\mathbf{0} &= -\mathbb{E}_{P(\mathbf{z}^t | \mathbf{y}^t)} (\mathbf{z}^t) + (1 + \kappa)\mu^t - \kappa(\nu^t - \mathbf{w}^t) \\
(1 + \kappa)\mu^t &= \mathbb{E}_{P(\mathbf{z}^t | \mathbf{y}^t)} (\mathbf{z}^t) + \kappa(\nu^t - \mathbf{w}^t) \\
\mu^t &= \frac{1}{1 + \kappa} \mathbb{E}_{P(\mathbf{z}^t | \mathbf{y}^t)} (\mathbf{z}^t) + \frac{\kappa}{1 + \kappa} (\nu^t - \mathbf{w}^t).
\end{aligned}$$

Evidently, the gradient of $\mathcal{L}(\phi, \mu)$ with respect to the graph decoder parameter ϕ is

$$\nabla_{\phi} \mathcal{L}(\phi, \mu) = - \sum_{t=1}^T \mathbb{E}_{P(\mathbf{z}^t | \mathbf{y}^t)} \left(\nabla_{\phi} \log P(\mathbf{y}^t | \mathbf{z}^t) \right).$$

432 7.2 Langevin Dynamics

Calculating the solution in (7) and the gradient in (8) requires evaluating the conditional expectations under the posterior distribution $P(\mathbf{z}^t | \mathbf{y}^t) \propto P(\mathbf{y}^t | \mathbf{z}^t) \times P(\mathbf{z}^t)$. In this section, we discuss the Langevin Dynamics

to sample $\mathbf{z}^t \in \mathbb{R}^d$ from the posterior distribution $P(\mathbf{z}^t|\mathbf{y}^t)$ that is conditional on the observed network $\mathbf{y}^t \in \{0, 1\}^{n \times n}$. The Langevin Dynamics, a short run MCMC, is achieved by iterating the following:

$$\begin{aligned} \mathbf{z}_{\tau+1}^t &= \mathbf{z}_\tau^t + \delta[\nabla_{\mathbf{z}^t} \log P(\mathbf{z}^t|\mathbf{y}^t)] + \sqrt{2\delta}\boldsymbol{\epsilon} \\ &= \mathbf{z}_\tau^t + \delta[\nabla_{\mathbf{z}^t} \log P(\mathbf{y}^t|\mathbf{z}^t) + \nabla_{\mathbf{z}^t} \log P(\mathbf{z}^t) - \nabla_{\mathbf{z}^t} \log P(\mathbf{y}^t)] + \sqrt{2\delta}\boldsymbol{\epsilon} \\ &= \mathbf{z}_\tau^t + \delta[\nabla_{\mathbf{z}^t} \log P(\mathbf{y}^t|\mathbf{z}^t) - (\mathbf{z}_\tau^t - \boldsymbol{\mu}^t)] + \sqrt{2\delta}\boldsymbol{\epsilon} \end{aligned}$$

433 where τ is the time step and δ is the step size of the Langevin Dynamics. The error term $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ serves
 434 as a random perturbation to the sampling process. The gradient of the graph decoder $P(\mathbf{y}^t|\mathbf{z}^t)$ with respect
 435 to the latent variable \mathbf{z}^t can be calculated efficiently through back-propagation. Essentially, we use MCMC
 436 samples to approximate the conditional expectation $\mathbb{E}_{P(\mathbf{z}^t|\mathbf{y}^t)}(\cdot)$ in the solution (7) and the gradient (8).

437 7.3 Group Lasso for Updating β

In this section, we present the derivation to update β in Proposition 2, which is equivalent to solving a Group Lasso problem Yuan & Lin (2006). We adapt the derivation from Bleakley & Vert (2011) for our proposed ADMM algorithm. Denote the objective function in (5) as $\mathcal{L}(\gamma, \beta)$. When $\beta_{t,\cdot} \neq \mathbf{0}$, the gradient of $\mathcal{L}(\gamma, \beta)$ with respect to $\beta_{t,\cdot}$ is

$$\nabla_{\beta_{t,\cdot}} \mathcal{L}(\gamma, \beta) = \lambda \frac{\beta_{t,\cdot}}{\|\beta_{t,\cdot}\|_2} - \kappa \mathbf{X}_{\cdot,t}^\top (\boldsymbol{\mu}_{(a+1)} + \mathbf{w}_{(a)} - \mathbf{1}_{T,1}\gamma - \mathbf{X}_{\cdot,t}\beta_{t,\cdot} - \mathbf{X}_{\cdot,-t}\beta_{-t,\cdot})$$

438 where $\mathbf{X}_{\cdot,t} \in \mathbb{R}^{T \times 1}$ is the t -th column of matrix $\mathbf{X} \in \mathbb{R}^{T \times (T-1)}$ and $\beta_{t,\cdot} \in \mathbb{R}^{1 \times d}$ is the t -th row of matrix
 439 $\beta \in \mathbb{R}^{(T-1) \times d}$. Moreover, we denote $\beta_{-t,\cdot} \in \mathbb{R}^{(T-1) \times p}$ as the matrix obtained by replacing the t -th row of
 440 matrix β with a zero vector, and $\mathbf{X}_{\cdot,-t} \in \mathbb{R}^{T \times (T-1)}$ is denoted similarly.

441 Setting the above gradient to zeros, we have

$$\beta_{t,\cdot} = (\kappa \mathbf{X}_{\cdot,t}^\top \mathbf{X}_{\cdot,t} + \frac{\lambda}{\|\beta_{t,\cdot}\|_2})^{-1} \mathbf{b}_t \quad (14)$$

where

$$\mathbf{b}_t = \kappa \mathbf{X}_{\cdot,t}^\top (\boldsymbol{\mu}_{(a+1)} + \mathbf{w}_{(a)} - \mathbf{1}_{T,1}\gamma - \mathbf{X}_{\cdot,-t}\beta_{-t,\cdot}) \in \mathbb{R}^{1 \times d}.$$

Calculating the Euclidean norm of (14) on both sides and rearrange the terms, we have

$$\|\beta_{t,\cdot}\|_2 = (\kappa \mathbf{X}_{\cdot,t}^\top \mathbf{X}_{\cdot,t})^{-1} (\|\mathbf{b}_t\|_2 - \lambda).$$

Plugging $\|\beta_{t,\cdot}\|_2$ into (14) for substitution, the solution of $\beta_{t,\cdot}$ is arrived at

$$\beta_{t,\cdot} = \frac{1}{\kappa \mathbf{X}_{\cdot,t}^\top \mathbf{X}_{\cdot,t}} \left(1 - \frac{\lambda}{\|\mathbf{b}_t\|_2}\right) \mathbf{b}_t.$$

Moreover, when $\beta_{t,\cdot} = \mathbf{0}$, the subgradient \mathbf{v} of $\|\beta_{t,\cdot}\|_2$ needs to satisfy that $\|\mathbf{v}\|_2 \leq 1$. Because

$$\mathbf{0} \in \lambda \mathbf{v} - \kappa \mathbf{X}_{\cdot,t}^\top (\boldsymbol{\mu}_{(a+1)} + \mathbf{w}_{(a)} - \mathbf{1}_{T,1}\gamma - \mathbf{X}_{\cdot,-t}\beta_{-t,\cdot}),$$

we obtain the condition that $\beta_{t,\cdot}$ becomes $\mathbf{0}$ when $\|\mathbf{b}_t\|_2 \leq \lambda$. Therefore, we can iteratively apply the following to update $\beta_{t,\cdot}$ for each block $t = 1, \dots, T-1$:

$$\beta_{t,\cdot} \leftarrow \frac{1}{\kappa \mathbf{X}_{\cdot,t}^\top \mathbf{X}_{\cdot,t}} \left(1 - \frac{\lambda}{\|\mathbf{b}_t\|_2}\right)_+ \mathbf{b}_t$$

442 where $(\cdot)_+ = \max(\cdot, 0)$.

443 7.4 Practical Guidelines

444 7.4.1 ADMM Implementation

445 In this section, we provide practical guidelines for the proposed framework and the Alternating Direction
 446 Method of Multipliers (ADMM) algorithm. For Langevin Dynamic sampling, we set $\delta = 0.5$, and we draw
 447 $s = 200$ samples for each time point t . To detect change points using the data-driven threshold in (13), we
 448 let the tuning parameter $\lambda = \{10, 20, 50, 100\}$. To detect change points using the localizing method with
 449 Gamma distribution in (11), we let the tuning parameter $\lambda = \{5, 10, 20, 50\}$. For each λ , we run $A = 50$
 450 iterations of ADMM. Within each ADMM iteration, we run $B = 20$ iterations of gradient descent with Adam
 451 optimizer for the graph decoder and $D = 20$ iterations of block coordinate descent for Group Lasso.

452 Since the proposed generative model is a probability distribution for the observed network data, in this work
 453 we stop ADMM learning with the following stopping criteria:

$$\left| \frac{l(\phi_{(a+1)}, \boldsymbol{\mu}_{(a+1)}) - l(\phi_{(a)}, \boldsymbol{\mu}_{(a)})}{l(\phi_{(a)}, \boldsymbol{\mu}_{(a)})} \right| \leq \epsilon_{\text{tol}}. \quad (15)$$

454 The log-likelihood $l(\phi, \boldsymbol{\mu})$ is approximated by sampling from the prior distribution $p(\mathbf{z}^t)$, as described in
 455 Section 4.2. Hence, we stop the ADMM procedure until the above criteria is satisfied for a' consecutive
 456 iterations. In Section 5, we set $\epsilon_{\text{tol}} = 10^{-5}$ and $a' = 5$.

Here we briefly elaborate on the computational aspect of the approximation of the log-likelihood. To calculate
 the product of edge probabilities for the conditional distribution $P(\mathbf{y}^t | \mathbf{z}^t)$, we have the following:

$$\begin{aligned} \sum_{t=1}^T \log P(\mathbf{y}^t) &= \sum_{t=1}^T \log \int P(\mathbf{y}^t | \mathbf{z}^t) P(\mathbf{z}^t) d\mathbf{z}^t \\ &= \sum_{t=1}^T \log \mathbb{E}_{P(\mathbf{z}^t)} \left[\prod_{(i,j) \in \mathbb{Y}} P(\mathbf{y}_{ij}^t | \mathbf{z}^t) \right] \\ &\approx \sum_{t=1}^T \log \left[\frac{1}{s} \sum_{u=1}^s \left[\prod_{(i,j) \in \mathbb{Y}} P(\mathbf{y}_{ij}^t | \mathbf{z}_u^t) \right] \right] \\ &= \sum_{t=1}^T \log \left[\frac{1}{s} \sum_{u=1}^s \exp \left\{ \sum_{(i,j) \in \mathbb{Y}} \log [P(\mathbf{y}_{ij}^t | \mathbf{z}_u^t)] \right\} \right] \\ &= \sum_{t=1}^T \left\{ -\log s + \log \left[\exp C^t \sum_{u=1}^s \exp \left\{ \sum_{(i,j) \in \mathbb{Y}} \log [P(\mathbf{y}_{ij}^t | \mathbf{z}_u^t)] - C^t \right\} \right] \right\} \\ &= \sum_{t=1}^T \left\{ C^t + \log \left[\sum_{u=1}^s \exp \left\{ \sum_{(i,j) \in \mathbb{Y}} \log [P(\mathbf{y}_{ij}^t | \mathbf{z}_u^t)] - C^t \right\} \right] \right\} - T \log s \end{aligned}$$

457 where $C^t \in \mathbb{R}$ is the maximum value of $\sum_{(i,j) \in \mathbb{Y}} \log [P(\mathbf{y}_{ij}^t | \mathbf{z}_u^t)]$ over m samples but within a time point t .

We also update the penalty parameter κ to improve convergence and to reduce reliance on its initialization.
 In particular, after the a -th ADMM iteration, we calculate the respective primal and dual residuals:

$$r_{\text{primal}}^{(a)} = \sqrt{\frac{1}{T \times d} \sum_{t=1}^T \|\boldsymbol{\mu}_{(a)}^t - \boldsymbol{\nu}_{(a)}^t\|_2^2} \quad \text{and} \quad r_{\text{dual}}^{(a)} = \sqrt{\frac{1}{T \times d} \sum_{t=1}^T \|\boldsymbol{\nu}_{(a)}^t - \boldsymbol{\nu}_{(a-1)}^t\|_2^2}.$$

Throughout, we initialize the penalty parameter $\kappa = 10$. We jointly update the penalty parameter κ and
 the scaled dual variable \mathbf{w} as in Boyd et al. (2011) with the following conditions:

$$\begin{aligned} \kappa_{(a+1)} &= 2\kappa_{(a)}, \quad \mathbf{w}_{(a+1)} = \frac{1}{2}\mathbf{w}_{(a)}, \quad \text{if } r_{\text{primal}}^{(a)} > 10 \times r_{\text{dual}}^{(a)}, \\ \kappa_{(a+1)} &= \frac{1}{2}\kappa_{(a)}, \quad \mathbf{w}_{(a+1)} = 2\mathbf{w}_{(a)}, \quad \text{if } r_{\text{dual}}^{(a)} > 10 \times r_{\text{primal}}^{(a)}. \end{aligned}$$

458 7.4.2 Post-Processing

Since neural networks may be over-fitted for a statistical model in change point detection, we track the following Coefficient of Variation as a signal-to-noise ratio when we learn the model parameter with the full data:

$$\text{Coefficient of Variation} = \frac{\text{mean}(\Delta\hat{\boldsymbol{\mu}})}{\text{sd}(\Delta\hat{\boldsymbol{\mu}})}.$$

459 We choose the learned parameter $\hat{\boldsymbol{\mu}}$ with the largest Coefficient of Variation as final output.

460 By convention, we also implement two post-processing steps to finalize the detected change points. When
461 the gap between two consecutive change points is small or $\hat{C}_k - \hat{C}_{k-1} < \epsilon_{\text{spc}}$, we preserve the detected change
462 point with greater $\Delta\hat{\zeta}$ value to prevent clusters of nearby change points. Moreover, as the endpoints of a
463 time span are usually not of interest, we remove the \hat{C}_k smaller than a threshold ϵ_{end} and the \hat{C}_k greater
464 than $T - \epsilon_{\text{end}}$. In Section 5, we set $\epsilon_{\text{spc}} = 5$ and $\epsilon_{\text{end}} = 5$.