Poisson-Gamma Dynamical Systems with Non-Stationary Transition Dynamics

Anonymous Author(s) Affiliation Address email

Abstract

Bayesian methodologies for handling count-valued time series have gained promi-1 2 nence due to their ability to infer interpretable latent structures and to estimate 3 uncertainties, and thus are especially suitable for dealing with *noisy* and *incomplete* count data. Among these Bayesian models, Poisson-Gamma Dynamical Systems 4 (PGDSs) are proven to be effective in capturing the evolving dynamics underlying 5 observed count sequences. However, the state-of-the-art PGDS still falls short in 6 capturing the *time-varying* transition dynamics that are commonly observed in 7 real-world count time series. To mitigate this limitation, a non-stationary PGDS 8 9 is proposed to allow the underlying transition matrices to evolve over time, and the evolving transition matrices are modeled by the specifically-designed Dirich-10 let Markov chains. Leveraging Dirichlet-Multinomial-Beta data augmentation 11 techniques, a fully-conjugate and efficient Gibbs sampler is developed to perform 12 posterior simulation. Experiments show that, in comparison with related models, 13 the proposed non-stationary PGDS achieves improved predictive performance 14 due to its capacity to learn non-stationary dependency structure captured by the 15 time-evolving transition matrices. 16

17 **1 Introduction**

In recent years, there has been an increasing interest in modeling count time series. For instance, 18 some previous works [1, 2, 3] are concerned with how to learn the evolving topics behind text 19 corpus (frequencies of words) over time. Some works [4, 5, 6, 7] try to predict global immigrant 20 trends underlying international population movements. Count time series are often overdispersed, 21 sparse, high-dimensional, and thus can not be well modeled by widely used dynamic models such 22 as linear dynamical systems [8, 9]. Recently, many works [10, 11, 12, 13, 14, 15, 16] prefer to 23 choose distributions of the gamma-Poisson family to build their hierarchical Bayesian models. In 24 particular, these models enjoy strong explainability and can estimate uncertainty especially when the 25 observations are *noisy* and *incomplete*. Among these works, Poisson-Gamma Dynamical Systems 26 (PGDSs) [13] received a lot of attention because PGDS can learn how the latent dimensions excite 27 each other to capture complicated dynamics in observed count series. For instance, a very inspiring 28 research paper may motivate other researchers to publish papers on related topics [17]. The outbreak 29 of COVID-19 in one state, may lead to the rapid rising of COVID-19 cases in the nearby states and 30 vice versa [18]. In particular, PGDS can be efficiently learned with a tractable Gibbs sampling scheme 31 via Poisson-Logarithmic data augmentation and marginalization technique [11]. Due to its strong 32 flexibility, PGDS achieves better performance in predicting missing entities and future observations, 33 compared with related models [9, 15]. 34

³⁵ Despite these advantages, PGDS still can not capture the time-varying transition dynamics underlying ³⁶ observed count sequences, which are commonly observed in real-world scenarios [19]. For instance,

Submitted to 38th Conference on Neural Information Processing Systems (NeurIPS 2024). Do not distribute.

during the initial stage of the COVID-19 pandemic, the worldwide counts of infectious patients were
 significantly affected by various local policies, government interventions, and emergent events [20,

³⁹ 21, 22]. The cross transition dynamics among the different monitoring areas were also evolving as

40 the corresponding policies and interventions changed over time. Hence, PGDS unavoidably makes a

41 certain amount of approximation error in capturing the aforementioned non-stationary count time

42 series, using a *time-invariant* transition kernel.

To mitigate this limitation, Non-Stationary Poisson-Gamma Dynamical Systems (NS-PGDSs), a novel
kind of Poisson-gamma dynamical systems with non-stationary transition dynamics are developed.
More specifically, NS-PGDS captures the evolving transition dynamics by the specifically-designed
Dirichlet Markov chains. Via the Dirichlet-Multinomial-Beta data augmentation strategy, the NonStationary Poisson-Gamma Dynamical Systems can be inferred with a conjugate-yet-efficient Gibbs
sampler. Our contributions are summarized as follows:

- We propose a Non-Stationary Poisson-Gamma Dynamical System (NS-PGDS), a novel
 Poisson-gamma dynamical system with time-evolving transition matrices that can well
 capture non-stationary transition dynamics underlying observed count series.
- Three Dirichlet Markov chains are dedicated to improving the flexibility and expressiveness of NS-PGDSs, for capturing the complex transition dynamics behind sequential count data.
- Fully-conjugate-yet-efficient Gibbs samplers are developed via Dirichlet-Multinomial-Beta
 augmentation techniques to perform posterior simulation for the proposed Dirichlet Markov
 chains.
- Extensive experiments are conducted on four real-world datasets, to evaluate the performance
 of the proposed NS-PGDS in predicting missing and future unseen observations. We also
 provide exploratory analysis to demonstrate the explainable latent structure inferred by the
 proposed NS-PGDS.

61 2 Preliminaries

Let $\boldsymbol{y}^{(t)} = \begin{bmatrix} y_1^{(t)}, \cdots, y_V^{(t)} \end{bmatrix}^{\mathrm{T}} \in \mathbb{N}^V$ be a vector of nonnegative count valued observations at time t. To capture the latent dynamics underlying count sequences, some previous works [23, 24] model the observations as

$$y^{(t)} = p(z^{(t)}), \ z^{(t)} = f^{-1}(x^{(t)}),$$

where $p(\cdot)$ is the observation likelihood function, and $f(\cdot)$ is an invertible link function that maps the parameters of observation component to continuous-valued latent variables $x^{(t)} \in \mathbb{R}^{K}$. The latent factor $x^{(t)}$ evolves over time according to a linear dynamical system (LDS) given by $x^{(t)} \sim$ $\mathcal{N}(Ax^{(t-1)}, \Lambda^{-1})$, where A is the state transition matrix of size $K \times K$, and $\Lambda = \text{diag}(\lambda_{1}, \cdots, \lambda_{K})$ is the inverse covariance matrix with λ_{k}^{-1} determining the variance of k-th latent dimension. Han et al. [23] adopted the Extended Rank likelihood function to model count observations using LDS with time complexity $\mathcal{O}((K + V)^{3})$, which prevents it from practical applications for analyzing large-scale count data.

Recently, Acharya et al. [15] and Schein et al. [13, 16] developed Poisson-gamma family models for sequential count observations. Gamma Process Dynamic Poisson Factor Analysis (GP-DPFA) [15] models count data as $y_v^{(t)} \sim \text{Pois}(\sum_{k=1}^K \lambda_k \phi_{vk} \theta_k^{(t)})$, where $\theta_k^{(t)}$ represents the strength of k-th latent factor at time t, and ϕ_{vk} captures the involvement degree of k-th factor to v-th observed dimension. To ensure the model identifiability, we can impose a restriction as $\sum_v \phi_{vk} = 1$, and thus place a Dirichlet prior over $\phi_k = [\phi_{1k}, \cdots, \phi_{Vk}]^T$ as $\phi_k \sim \text{Dir}(\epsilon_0, \cdots, \epsilon_0)$.

Markov chain as $\theta_k^{(t)} \sim \text{Gam}(\theta_k^{(t-1)}, c_t)$, where c_t is the rate parameter of the gamma distribution to control the variance of the gamma Markov chains. Although GP-DPFA can well fit one-dimensional count sequences, it fails to learn how the latent dimensions interact with each other.

⁸³ To address this concern, Schein et al. [13] developed Poisson-gamma dynamical systems to ⁸⁴ capture the underlying transition dynamics. In particular, $\theta_k^{(t)}$ evolves over time as $\theta_k^{(t)} \sim$



Figure 1: The graphical representation of the NS-PGDS. The time interval is divided into equallyspaced sub-intervals. Each sub-interval contains M time steps. The transition dynamics is stationary within a sub-interval. In particular, the transition matrices evolve over sub-intervals via Dirichlet Markov processes while latent factors evolve over time steps via Eq.(1).

Gam $(\tau_0 \sum_{k_2=1}^K \pi_{kk_2} \theta_{k_2}^{(t-1)}, \tau_0)$, where π_{kk_2} represents how k_2 -th latent factor excites the k-th latent factor at next time step, and $\sum_{k=1}^K \pi_{kk_2} = 1$.



3 Non-Stationary Poisson-Gamma Dynamical Systems

Figure 2: An example illustrates the Poisson-gamma dynamical systems with non-stationary transition kernels. The three gamma dynamic processes independently evolve over time during the (i - 1)-th interval. During *i*-th interval, $\theta_1^{(t)}$ and $\theta_2^{(t)}$ gradually starts to interact with each other while $\theta_3^{(t)}$ remains independent to the other two dimensions. During (i + 1)-th interval all the three latent components

(i + 1) in merval an the three fatences start to interact with each other.

103 104 transition dynamics. For instance, the transition dynamics behind COVID-19 infectious processes are time-varying, and highly affected by various interventional policies. Hence, to mitigate this limitation, we model the count sequences as $y_v^{(t)} \sim \text{Pois}\left(\delta^{(t)} \sum_{k=1}^K \phi_{vk} \theta_k^{(t)}\right),$

Real-world count time sequences are often non-

stationary because the external interventional

environments are always changing over time.

The stationary PGDS with a time-invariant tran-

sition kernel fails to capture such time-varying

in which, the latent factors are specified by

$$\theta_k^{(t)} \sim \operatorname{Gam}\left(\tau_0 \sum_{k_2=1}^K \pi_{kk_2}^{(t-1)} \theta_{k_2}^{(t-1)}, \tau_0\right),\tag{1}$$

where the multiplicative term $\delta^{(t)} \sim \text{Gam}(\epsilon_0, \epsilon_0)$ and the transition matrices are time-varying as $\mathbf{\Pi}^{(t)} \equiv \left[\pi_{kk_2}^{(t)}\right]_{k,k_2=1}^{K}$. As shown in Figure 2, to model the time-varying transition dynamics, we assume the whole time

interval can be divided into *I* equally-spaced sub-intervals. The transition kernel behind complicated dynamic counts is assumed to be *static* within each sub-interval, while evolving over sub-intervals, to capture non-stationary behaviours. In another word, the proposed model allows the latent factors to evolve over time steps while the transition matrices change over sub-intervals but assumed to be stationary within each sub-interval, as shown in Figure 1. In particular, we let each sub-interval contains *M* time steps, and the *i*-th interval contains time steps $\{t \mid t = (i-1)M + 1, \dots, iM\}$. We define i(t) as the function that maps time step t to its corresponding sub-interval.

Dirichlet-Dirichlet Markov processes. To capture how the underlying transition kernel smoothly
 evolves over sub-intervals, we first propose the Dirichlet-Dirichlet (Dir-Dir) Markov chain as

$$\boldsymbol{\pi}_{k}^{(i)} \mid \boldsymbol{\pi}_{k}^{(i-1)} \sim \operatorname{Dir}\left(\eta K \boldsymbol{\pi}_{1k}^{(i-1)}, \cdots, \eta K \boldsymbol{\pi}_{Kk}^{(i-1)}\right),$$
(2)

where $\pi_k^{(i)}$ represents the k-th column of $\Pi^{(i)}$, and the prior of the scaling parameter η is given by $\eta \sim \text{Gam}(e_0, f_0).$ The initial states are defined as $\theta_k^{(1)} \sim \text{Gam}(\tau_0\nu_k,\tau_0)$. The prior for the transition kernel of the first sub-interval is given by $\pi_k^{(1)} \sim \text{Dir}(\nu_1\nu_k,\cdots,\xi\nu_k,\cdots,\nu_K\nu_k)$, where $\nu_k \sim \text{Gam}(\frac{\gamma_0}{K},\beta)$ and $\xi,\beta \sim \text{Gam}(\epsilon_0,\epsilon_0)$. Note that the expectation and variance of the transition kernel at *i*-th sub-interval can be calculated as

$$\mathsf{E}\left[\pi_{k}^{(i)} \mid \pi_{k}^{(i-1)}\right] = \pi_{k}^{(i-1)}, \quad \mathsf{Var}\left[\pi_{k_{1}k}^{(i)} \mid \pi_{k}^{(i-1)}\right] = \frac{\pi_{k_{1}k}^{(i-1)}\left(1 - \pi_{k_{1}k}^{(i-1)}\right)}{\eta K + 1},$$

respectively. The transition dynamics of *i*-th sub-interval inherits the information of the previous sub-interval, and also adapts to the data observed in the current sub-interval. The scaling parameter η controls the variance of the transition matrices.

The prior specification defined in Eq.(2) by rescaling the transition matrix at the previous sub-interval allows the transition dynamics to change smoothly, and thus might be insufficient to capture the rapid changes observed in complicated dynamics. To further improve the flexibility of the transition structure, two modified Dirichlet Markov chains are studied to capture the correlation structure between the dimensions of the transition matrices over time. **Dirichlet-Gamma-Dirichlet Markov processes.** We first

introduce the Dirichlet-Gamma-Dirichlet (Dir-Gam-Dir)
 Markov chain to model the evolving transition matrices
 as

$$\pi_{k}^{(i)} \sim \operatorname{Dir}\left(\alpha_{1k}^{(i)}, \cdots, \alpha_{Kk}^{(i)}\right),$$

$$\alpha_{k_{1}k}^{(i)} \sim \operatorname{Gam}\left(\gamma_{k}^{(i-1)} \sum_{k_{2}=1}^{K} \psi_{kk_{1}k_{2}}^{(i-1)} \pi_{k_{2}k}^{(i-1)}, c_{k}^{(i)}\right), \quad (3)$$

where we use $\psi_{kk_1k_2}^{(i-1)}$ to capture the mutation between two consecutive sub-intervals, and its prior is given by

$$\left(\psi_{k1k_2}^{(i-1)}, \cdots, \psi_{kKk_2}^{(i-1)}\right) \sim \operatorname{Dir}\left(\epsilon_0, \cdots, \epsilon_0\right),$$

and $\gamma_k^{(i)}, c_k^{(i)} \sim \text{Gam}(\epsilon_0, \epsilon_0)$. Compared with the construction defined by Eq.(2), the expectation of Dirichlet-Gamma-Dirichlet Markov chain is

$$\mathsf{E}\left[\boldsymbol{\pi}_{k}^{(i)} \mid \boldsymbol{\pi}_{k}^{(i-1)}\right] = \boldsymbol{\Psi}_{k}^{(i-1)} \boldsymbol{\pi}_{k}^{(i-1)}.$$



(: 1) (

(: 1))

Figure 3: Diagrams of the proposed Dirichlet Markov constructions. (a) is the Dir-Dir construction. (b) is the Dir-Gam-Dir construction which takes mutation into account. (c) illustrates the PR-Gam-Dir construction which adopts Poisson randomized gamma distribution and can be equivalently represented as Eq.(5).

This construction takes interactions among components of columns into account. Hence it will dramatically improve the flexibility of our model and thus better fit more complicated dynamics, compared with Dir-Dir Markov chains that only yield smoothing transition dynamics.

Poisson-randomized-gamma-Dirichlet Markov processes. By leveraging the Poisson-randomized
 gamma distribution [25], we introduce another type of time-varying transition kernels, which also
 model the interactions among components like Dir-Gam-Dir construction but may induce different
 properties such as sparsity. The Poisson-randomized-gamma-Dirichlet (PR-Gam-Dir) Markov chain
 can be formulated as

$$\boldsymbol{\pi}_{k}^{(i)} \sim \operatorname{Dir}\left(\alpha_{1k}^{(i)}, \cdots, \alpha_{Kk}^{(i)}\right), \ \alpha_{k_{1}k}^{(i)} \sim \operatorname{RG1}\left(\boldsymbol{\epsilon}^{\alpha}, \gamma_{k}^{(i-1)} \sum_{k_{2}=1}^{K} \psi_{kk_{1}k_{2}}^{(i-1)} \boldsymbol{\pi}_{k_{2}k}^{(i-1)}, \boldsymbol{c}_{k}^{(i)}\right),$$
(4)

where RG1 (·) denotes the randomized gamma distribution of the first type. Similarly, for $\psi_{kk_1k_2}^{(i-1)}$, $\gamma_k^{(i)}$, and $c_k^{(i)}$, the priors are given by

$$\left(\psi_{k1k_2}^{(i-1)}, \cdots, \psi_{kKk_2}^{(i-1)}\right) \sim \operatorname{Dir}\left(\epsilon_0, \cdots, \epsilon_0\right), \ \gamma_k^{(i)}, c_k^{(i)} \sim \operatorname{Gam}\left(\epsilon_0, \epsilon_0\right), \ \operatorname{respectively.}$$

¹⁴⁷ The diagrams of three Dirichlet Markov constructions are shown in Figure 3.

4 Markov Chain Monte Carlo Inference

¹⁴⁹ In this section, we present the Gibbs sampler for the proposed NS-PGDS. We only illustrate the key ¹⁵⁰ points of the derivation and the details can be found in the appendix.

- 151 **Lemma 1** If $y \sim \text{NB}(a, g(\zeta))$ and $l \sim \text{CRT}(y, a)$, where $\text{NB}(\cdot)$ refers to negative-binomial
- *distribution*, CRT (·) *represents Chinese restaurant table distribution* [26], and $g(z) = 1 \exp(-z)$.
- 153 Then the joint distribution of y and l can be equivalently distributed as $y \sim \text{SumLog}(l, g(\zeta))$ and
- 154 $l \sim \text{Pois}(a\zeta)$ [11], *i.e.*

 $NB(y; a, g(\zeta)) CRT(l; y, a) = SumLog(y; l, g(\zeta)) Pois(l; a\zeta),$

where SumLog $(l, g(\zeta)) = \sum_{i=1}^{l} x_i$ and $x_i \sim \text{Log}(g(\zeta))$ are independently and identically logarithmic distributed random variables [27].

Lemma 2 Suppose $\mathbf{n} = (n_1, \cdots, n_K)$ and

$$\mathbf{n} \mid n \sim \text{DirMult}(n, r_1, \cdots, r_K),$$

- where $DirMult(\cdot)$ refers to Dirichlet-multimonial distribution. We sample the augmented variable
- 158 $q \mid n \sim \text{Beta}(n, r)$, where $r = \sum_{k=1}^{K} r_k$. According to [28], conditioning on q, we have 159 $n_k \sim \text{NB}(r_k, q)$.
- Sampling $y_{vk}^{(t)}$: Use the relationship between Poisson and multinomial distributions, we sample

$$\left(\left(y_{vk}^{(t)} \right)_{k=1}^{K} \mid - \right) \sim \operatorname{Mult} \left(y_{v}^{(t)}, \left(\frac{\phi_{vk} \theta_{k}^{(t)}}{\sum_{k=1}^{K} \phi_{vk} \theta_{k}^{(t)}} \right)_{k=1}^{K} \right)$$

161 Sampling ϕ_k : Via Dirichlet-multinomial conjugacy, the posterior of ϕ_k is

$$(\boldsymbol{\phi}_{\boldsymbol{k}} \mid -) \sim \operatorname{Dir}\left(\epsilon_0 + \sum_{t=1}^T y_{1k}^{(t)}, \cdots, \epsilon_0 + \sum_{t=1}^T y_{Vk}^{(t)}\right)$$

162 **Sampling** $\theta_k^{(t)}$: To sample from the posterior of $\theta_k^{(t)}$, we first sample the auxiliary variables. Setting 163 $l_{\cdot k}^{(T+1)} = 0$ and $\zeta^{(T+1)} = 0$, we sample the augmented variables backwards from $t = T, \dots, 2$,

$$\begin{pmatrix} l_{k}^{(t)} \mid - \end{pmatrix} \sim \operatorname{CRT} \left(y_{\cdot k}^{(t)} + l_{\cdot k}^{(t+1)}, \tau_0 \sum_{k_2=1}^{K} \pi_{k k_2}^{i(t-1)} \theta_{k_2}^{(t-1)} \right), \\ \begin{pmatrix} l_{k_1}^{(t)}, \cdots, l_{kK}^{(t)} \mid - \end{pmatrix} \sim \operatorname{Mult} \left(l_{k}^{(t)}, \left(\frac{\pi_{k_1}^{i(t-1)} \theta_1^{(t-1)}}{\sum_{k_2=1}^{K} \pi_{k k_2}^{i(t-1)} \theta_{k_2}^{(t-1)}}, \cdots, \frac{\pi_{kK}^{i(t-1)} \theta_K^{(t-1)}}{\sum_{k_2=1}^{K} \pi_{k k_2}^{i(t-1)} \theta_{k_2}^{(t-1)}} \right) \right).$$

Let us define $l_{k}^{(t)} = \sum_{k_{1}=1}^{K} l_{k_{1}k}^{(t)}$ and $\zeta^{(t)} = \ln(1 + \frac{\delta^{(t)}}{\tau_{0}} + \zeta^{(t+1)})$. After sampling the auxiliary variables, then for $t = 1, \dots, T$, by Poisson-gamma conjugacy, we obtain

$$\left(\theta_{k}^{(t)} \mid -\right) \sim \operatorname{Gam}\left(y_{\cdot k}^{(t)} + l_{\cdot k}^{(t+1)} + \tau_0 \sum_{k_2=1}^{K} \pi_{k k_2}^{i(t-1)} \theta_{k_2}^{(t-1)}, \tau_0 + \delta^{(t)} + \zeta^{(t+1)} \tau_0\right)$$

Sampling $\Pi^{(i)}$: We only illustrate Gibbs sampling algorithm for PR-Gam-Dir construction, sampling algorithms for other constructions can be found in the appendix. We define M as the length of each sub-interval, and I as the number of intervals. For $i = I, \dots, 2$, because $(l_{1k}^{(i)}, \dots, l_{Kk}^{(i)})$ and $(g_{1k}^{(i+1)}, \dots, g_{KK}^{(i+1)})$ are multinomially distributed, where $l_{k_1k}^{(i)} = \sum_{(i-1)M+1}^{iM} l_{k_1k}^{(i)}$ refers to the summation of $l_{k_1k}^{(i)}$ over *i*-th sub-interval and same notation for other variables. By the definition of Dirichlet-multinomial distribution and Lemma 2, defining $g_{k_1k}^{(I+1)} = 0$, we sample the auxiliary variables as $(q_k^{(i)} | -) \sim \text{Beta}(l_{\cdot k}^{(i)} + g_{\cdot k}^{(i+1)}, \alpha_{\cdot k}^{(i)})$, then we have $(l_{k_1k}^{(i)} + g_{\cdot k_1k}^{(i+1)}) \sim \text{NB}(\alpha_{k_1k}^{(i)}, q_k^{(i)})$. Then we further sample $(h_{k_1k}^{(i)} | -) \sim \text{CRT}(l_{k_1k}^{(i)} + g_{\cdot k_1k}^{(i+1)}, \alpha_{k_1k}^{(i)})$. Via Lemma 1, we obtain $h_{k_1k}^{(i)} \sim$ Pois $(-\alpha_{k_1k}^{(i)} \ln(1 - q_k^{(i)}))$. For Dirichlet-Randomized-Gamma-Dirichlet Markov construction defined by Eq.(4), we can equivalently represent it as

$$\alpha_{k_1k}^{(i)} \sim \operatorname{Gam}\left(g_{k_1k}^{(i)} + \epsilon^{\alpha}, c_k^{(i)}\right), \ g_{k_1k}^{(i)} = \operatorname{Pois}\left(\gamma^{(i-1)} \sum_{k_2=1}^{K} \psi_{kk_1k_2}^{(i-1)} \pi_{k_2k}^{(i-1)}\right).$$
(5)

We define $\lambda_{k_1k}^{(i-1)} \triangleq \gamma_k^{(i-1)} \sum_{k_2=1}^K \psi_{kk_1k_2}^{(i-1)} \pi_{k_2k}^{(i-1)}$ for notation conciseness. By Poisson-gamma conjugacy, we have $(\alpha_{k_1k}^{(i)} \mid -) \sim \text{Gam}(g_{k_1k}^{(i)} + \epsilon^{\alpha} + h_{k_1k}^{(i)}, c_k^{(i)} - \ln(1 - q_k^{(i)}))$. If $\epsilon^{\alpha} > 0$, we can

sample the posterior of $g_{k_1k}^{(i)}$ via $(g_{k_1k}^{(i)} | -) \sim \text{Bessel}(\epsilon^{\alpha} - 1, 2\sqrt{\alpha_{k_1k}^{(i)}c_k^{(i)}\lambda_{k_1k}^{(i-1)}})$, where $\text{Bessel}(\cdot)$ denotes Bessel distribution. If $\epsilon^{\alpha} = 0$, we sample $g_{k_1k}^{(i)}$ via

$$\begin{pmatrix} g_{k_1k}^{(i)} \mid - \end{pmatrix} \sim \begin{cases} & \operatorname{Pois}\left(\frac{c_k^{(i)}\lambda_{k_1k}^{(i-1)}}{c_k^{(i)} - \ln\left(1 - q_k^{(i)}\right)}\right) & \text{if } h_{k_1k}^{(i)} = 0 \\ & \\ & \operatorname{SCH}\left(h_{k_1k}^{(i)}, \frac{c_k^{(i)}\lambda_{k_1k}^{(i-1)}}{c_k^{(i)} - \ln\left(1 - q_k^{(i)}\right)}\right) & \text{otherwise,} \end{cases}$$

where SCH (·) denotes the shifted confluent hypergeometric distribution [16]. Defining $g_{k_1k}^{(i)} = g_{k_1\cdot k}^{(i)} = \sum_{k2=1}^{K} g_{k_1k_2k}^{(i)}$, we first augment

$$\left(g_{k_{1}1k}^{(i)}, \cdots, g_{k_{1}Kk}^{(i)}\right) \sim \operatorname{Mult}\left(g_{k_{1}k}^{(i)}, \left(\psi_{kk_{1}k_{2}}^{(i-1)} \pi_{k_{2}k}^{(i-1)}\right)_{k_{2}=1}^{K}\right)$$

then we obtain $g_{k_1k_2k}^{(i)} \sim \text{Pois}(\gamma^{(i-1)}\psi_{kk_1k_2}^{(i-1)}\pi_{k_2k}^{(i-1)})$. By Dirichlet-multinomial conjugacy, we have

$$\left(\left(\psi_{k1k_2}^{(i-1)}, \cdots, \psi_{kKk_2}^{(i-1)} \right) \mid - \right) \sim \operatorname{Dir} \left(\epsilon_0 + g_{1k_2k}^{(i)}, \cdots, \epsilon_0 + g_{Kk_2k}^{(i)} \right), \text{ and} \left(\pi_k^{(i-1)} \mid - \right) \sim \operatorname{Dir} \left(\alpha_{1k}^{(i-1)} + l_{1k}^{(i-1)} + g_{\cdot 1k}^{(i)}, \cdots, \alpha_{Kk}^{(i-1)} + l_{Kk}^{(i-1)} + g_{\cdot Kk}^{(i)} \right).$$

183 Specifically, we have $\alpha_{k_1k}^{(1)} = \nu_{k_1}\nu_k$, if $k_1 \neq k$, and $\alpha_{k_1k}^{(1)} = \xi\nu_k$, if $k_1 = k$.

184 **5 Related Work**

Modeling count time sequences has been receiving increasing attentions in statistical and machine 185 learning communities. Han et al. [23] adopted linear dynamical systems to capture the underlying 186 dynamics of the data and leveraged Extended Rank likelihood function to model count observations. 187 Some Poisson-gamma models assume that the count vector at each time step is modeled by Poisson 188 factor analysis (PFA) [11] and leverage special stochastic processes to model the temporal dependen-189 cies of latent factors. For example, gamma process dynamic Poisson factor analysis (GP-DPFA) [15] 190 adopts gamma Markov chains which assumes the latent factor of the next time step is drawn from 191 a gamma distribution with the shape parameter be the latent factor of the current time step. Schein 192 et al. [13] proposed Poisson-gamma dynamical systems (PGDSs), which take the interactions among 193 latent dimensions into account and use a transition matrix to capture the interactions. Deep dynamic 194 Poisson factor analysis (DDPFA) [29] adopts recurrent neural networks (RNNs) to capture the com-195 plex long-term dependencies of latent factors. Yang and Koeppl [30] applied Poisson-gamma count 196 model to analyze relational data arising from longitudinal networks, which can capture the evolution 197 of individual node-group memberships over time. Many modifications of PGDS have been proposed 198 in recent years. Guo et al. [31] proposed deep Poisson-gamma dynamical systems which aim to 199 capture the long-range temporal dependencies. Schein et al. [16] employed Poisson-randomized 200 gamma distribution to build a new transition process of latent factors. Chen et al. [32] proposed 201 Switching Poisson-gamma dynamical systems (SPGDS), allowing PGDS to select from several tran-202 sition matrices, and thus can better adapt to nonlinear dynamics. In contrast to SPGDS, the number 203 of transition matrices of the proposed NS-PGDS is not limited and thus can be adopted to analyze 204 various complicated non-stationary count sequences. Filstroff et al. [33] extensively analyzed many 205 gamma Markov chains for non-negative matrix factorization and introduced new gamma Markov 206 chains with well-defined stationary distribution (BGAR). 207

208 6 Experiments

We conducted experiments for both predictive and exploratory analysis to demonstrate the ability of the proposed model in capturing non-stationary count time sequences. The baseline models included in the experiments are: 1) Gamma process dynamic Poisson factor analysis (GP-DPFA) [15]. GP-DPFA models the evolution of latent components as $\theta_k^{(t)} \sim \text{Gam}(\theta_k^{(t-1)}, c_t)$, in which each component evolves independently of the other components. 2) Gamma Markov chains on the rate parameter of gamma distribution (GMC-RATE) [33]. GMC-RATE adopts gamma Markov chains defined via the rate parameter of the gamma distribution to model the evolution of $\theta_k^{(t)}$

			GP-DPFA	GMC-RATE	GMC-HIER	BGAR	PGDS	NS-PGDS (Dir-Dir)	NS-PGDS (Dir-Gam-Dir)	NS-PGDS (PR-Gam-Dir)
ICEWS	MAE	S	0.259 ± 0.005	0.258 ± 0.005	0.256 ± 0.006	0.264 ± 0.006	0.215 ± 0.007	0.215 ± 0.008	$\textbf{0.214} \pm 0.008$	0.215 ± 0.008
		F	0.176 ± 0.005	0.187 ± 0.003	0.185 ± 0.016	0.222 ± 0.043	0.185 ± 0.003	0.167 ± 0.009	0.169 ± 0.006	0.169 ± 0.009
	MRE	S	0.125 ± 0.003	0.124 ± 0.002	0.122 ± 0.003	0.130 ± 0.004	0.102 ± 0.005	0.101 ± 0.005	0.101 ± 0.005	0.102 ± 0.005
		F	0.099 ± 0.006	0.114 ± 0.003	0.111 ± 0.018	0.142 ± 0.036	0.108 ± 0.001	0.094 ± 0.005	0.097 ± 0.004	0.097 ± 0.008
NIPS	MAE	S	$18.299 \pm \! 6.545$	$17.105 \ {\pm 6.449}$	17.098 ± 6.441	$17.935 \ {\pm 6.450}$	$14.706 \pm \scriptscriptstyle 4.414$	$14.032 \pm \!$	$14.026 \pm \scriptscriptstyle 4.405$	$14.014 \pm \scriptscriptstyle 4.387$
		F	$48.355 \ {\pm 1.461}$	$46.234 \ {\scriptstyle \pm 1.629}$	102.506 ± 39.932	$62.449 \ {\scriptstyle \pm 14.463}$	51.562 ± 0.679	45.979 ± 1.342	46.710 ± 1.152	$46.582 \ \pm 1.196$
	MRE	S	0.729 ± 0.412	0.684 ± 0.316	0.664 ± 0.315	0.769 ± 0.366	0.590 ± 0.097	0.581 ± 0.090	0.581 ± 0.090	0.580 ±0.090
		F	0.415 ± 0.016	0.387 ± 0.023	0.580 ± 0.148	0.465 ± 0.049	0.459 ± 0.006	0.399 ± 0.003	0.395 ± 0.006	0.397 ± 0.003
USEI	MAE	S	4.681 ± 0.564	$4.931{\scriptstyle~\pm 0.872}$	4.748 ± 0.829	5.244 ± 0.939	4.703 ± 0.538	$4.600{\scriptstyle~\pm 0.542}$	4.608 ± 0.541	4.596 ± 0.562
		F	11.665 ± 0.367	9.454 ± 0.809	12.423 ± 1.060	21.948 ± 0.133	11.118 ± 0.220	7.973 ± 1.222	7.168 ±1.221	7.296 ± 1.127
	MRE	S	1.458 ± 0.177	1.128 ± 0.189	1.088 ± 0.162	1.941 ± 0.209	1.279 ± 0.257	1.309 ± 0.220	1.298 ± 0.236	1.301 ± 0.229
		F	7.473 ± 0.623	6.508 ± 0.571	8.929 ± 2.514	13.706 ± 1.268	4.238 ± 0.325	2.602 ± 0.455	2.577 ± 0.331	2.685 ± 0.366
COVID-19	MAE	S	7.935 ± 0.751	7.144 ± 1.159	$7.240{\scriptstyle~\pm 0.848}$	$7.819{\scriptstyle~\pm1.348}$	$7.566{\scriptstyle~\pm1.095}$	6.969 ±1.107	$6.988 {\scriptstyle~\pm 1.056}$	$6.981{\scriptstyle~\pm1.022}$
		F	9.137 ± 1.102	9.600 ± 1.257	10.409 ± 1.910	12.550 ± 2.156	9.314 ± 0.236	8.799 ± 0.706	8.770 ±0.438	9.033 ± 0.477
	MRE	S	0.564 ± 0.126	0.493 ± 0.136	0.504 ± 0.109	0.769 ± 0.169	0.558 ± 0.130	0.523 ± 0.125	0.525 ± 0.124	0.526 ± 0.123
		F	0.627 ± 0.106	0.556 ± 0.052	0.585 ± 0.067	0.759 ± 0.150	0.585 ± 0.007	0.523 ± 0.028	$0.519{\scriptstyle~\pm 0.017}$	$\textbf{0.513} \pm 0.014$

Table 1: Results of predictive analysis. "S" means data smoothing and "F" means data forecasting.

as $\theta_k^{(t)} \sim \text{Gam}(\alpha, \beta/\theta_k^{(t-1)})$. 3) Gamma Markov chains on the rate parameter with hierarchical auxiliary variable (GMC-HIER) [33]. GMC-HIER models the evolution of latent components with an auxiliary variables as $z_k^{(t)} \sim \text{Gam}(\alpha_z, \beta_z \theta_k^{(t-1)})$ and $\theta_k^{(t)} \sim \text{Gam}(a_\theta, \beta_\theta z_k^{(t)})$. 4) Autogressive beta-gamma process (BGAR) [34, 33]. BGAR is also a gamma Markov model. In contrast to the above models, there is a well-defined stationary distribution for BGAR. 5) Poisson-gamma dynamical system (PGDS) [13] takes interactions among latent dimensions into account, and models the evolution of $\theta_k^{(t)}$ as $\theta_k^{(t)} \sim \text{Gam}(\tau_0 \sum_{k_2=1}^K \pi_{kk_2} \theta_{k_2}^{(t-1)}, \tau_0)$.

The real-world datasets used in the experiments are: 1) Integrated Crisis Early Warning System 223 (ICEWS): ICEWS is an international relations event dataset, comprising interaction events between 224 countries extracted from news corpora. For ICEWS dataset, we have T = 365 time steps and 225 V = 6197 dimensions, and we set M = 30. 2) **NIPS**: NIPS dataset contains the papers published in 226 the NeurIPS conference from 1987 to 2015. We have T = 28 time steps and V = 2000 dimensions 227 for NIPS dataset and we set M = 5. 3) U.S. Earthquake Intensity (USEI): USEI contains a 228 collection of damage and felt reports for U.S. (and a few other countries) earthquakes. We use the 229 monthly reports from 1957-1986 and have T = 348, V = 64 and set M = 34. 4) **COVID-19**: This 230 dataset contains daily death cases data for states in the United States, spanning from March 2020 to 231 June 2020. For this dataset, we have V = 51 dimensions and T = 90 time steps and set M = 20. 232

233 6.1 Predictive Analysis

To compare the predictive performance of the proposed model with the baselines, we considered two 234 standard tasks: data smoothing and forecasting. For data smoothing task, our objective is to predict 235 $y^{(t)}$ given the remaining data observation $Y \setminus y^{(t)}$. To this end, we randomly masked 10 percents of 236 the observed data over non-adjacent time steps, and predicted the masked values. For forecasting task, we held out data of the last S time steps, and predicted $y^{(T+1)}, \dots, y^{(T+S)}$ given $y^{(1)}, \dots, y^{(T)}$. In 237 238 this experiment we set S = 2. We ran the baseline models including GP-DPFA, PGDS, GMC-RATE, 239 GMC-HIER, BGAR, using their default settings as provided in [15, 13, 33]. For the NS-PGDS, we set 240 K = 100 for ICEWS, K = 10 for other datasets, and set $\tau_0 = 1, \gamma_0 = 50, \epsilon_0 = 0.1$. We performed 241 4000 Gibbs sampling iterations. In the experiments, we found that the Gibbs sampler started to 242 converge after 1000 iterations, and thus we set the burn-in time be 2000 iterations. We retained 243 244 every hundredth sample, and averaged the predictions over the samples. Mean relative error (MRE) and mean absolute error (MAE) are adopted to evaluate the model's predictive capability, which are defined as $MRE = \frac{1}{TV} \sum_{t} \sum_{v} \frac{|y_v^{(t)} - \hat{y}_v^{(t)}|}{1 + y_v^{(t)}}$ and $MAE = \frac{1}{TV} \sum_{t} \sum_{v} |y_v^{(t)} - \hat{y}_v^{(t)}|$ respectively, where $y_v^{(t)}$ indicates the true count and $\hat{y}_v^{(t)}$ is the prediction. 245 246 247

As the experiment results shown in Table 1, the NS-PGDS exhibits improved performance in both data smoothing and forecasting tasks. We attribute this enhanced capability to the time-varying transition kernels, which effectively adapt to the non-stationary environment, and thus achieve improved predictive performance. For some datasets (e.g. ICEWS) and tasks, the effectiveness of the Dir-Gam-Dir and Pr-Gam-Dir constructions does not be exhibited in the numerical results. However, these two constructions indeed induce more informative patterns compared with Dir-Dir construction, as shown in the exploratory analysis.



Figure 4: The latent factors inferred by the NS-PGDS. (a) and (b) illustrate the top 2 latent factors inferred from ICEWS dataset, (a) corresponds to Iraq war and (b) corresponds to the Six-Party Talks. (c) illustrates the evolving trends of the top 5 latent factors inferred from NIPS dataset.

We used ICEWS and NIPS datasets for exploratory analysis, and chose the NS-PGDS with Dirichlet-256 Dirichlet Markov chains for illustration. Figure 4(a) and Figure 4(b) demonstrate the top 2 latent 257 factors inferred by NS-PGDS from ICEWS dataset. From Figure 4(a) we can see that the main labels 258 are "Iraq (IRQ)-United States (USA)", "Iraq (IRQ)-United Kingdom (UK)", "Russia (RUS)-United 259 States (USA)", and so on. This latent factor probably corresponds to the topic about Iraq war. Besides, 260 261 in Figure 4(a), there is a peak around March, 2003, and we know that the Iraq war broke out exactly on 20 March, 2003. In addition, the most dominant labels shown in Figure 4(b) are "Japan (JPN)–United 262 States (USA)", "China (CHN)–United States (USA)", "North Korea (PRK)–United States (USA)", 263 "South Korea (KOR)–United States (USA)", and so on. We can infer that this latent factor corresponds 264 to "Six-Party Talks" and other accidents about it. 265

Figure 4(c) demonstrates the evolving trends of the top 5 latent factors inferred by the NS-PGDS from NIPS dataset, and the legend indicates the representative words of the corresponding latent factors. Clearly, the green and blue lines correspond to the latent factors of neural network research which started to decline from the 1990s. From the 1990s we see that the latent factors about statistical and probabilistic methods began to dominate the NeurIPS conference. In addition, the NS-PGDS also captured the revival of neural networks (blue line) from the 2010s. The above observations from the latent structure inferred by the NS-PGDS match our prior knowledge.



Figure 5: Transition matrices inferred from NIPS dataset. (a)
illustrates the transition matrix inferred by the PGDS. (b)-(f)
illustrate the time-varying transition matrices inferred by the
NS-PGDS.

292

Next, we explored the time-varying transition matrices inferred by the NS-PGDS. We chose NIPS dataset for illustratiuon, and set K = 10 and the interval length M to be 5. The time-varying transition matrices are shown from Figure 5(b) to Figure 5(f). At the beginning, matrices shown in Figure 5(b) and Figure 5(c) are close to identity matrices. Then the transition matrices tend to become block diagonal matrices with 2 blocks, as shown in Figure 5(d)-5(f). The representative words for latent factors in the first block "state-linear-classification", are "network-neural-networks", "kernelimage-space", "network-neuralnetworks", "neural-networks-state".

The representative words for latent factors in the second block are "image-sparse-matrix", "kernelsupervised-random", "matrix-sample-random", "inference-prior-latent", "state-policy-gamma". The first block primarily captured the correlations among the research topics about neural networks. The second block reflects that, from the 1990s, statistical learning and Bayesian methods began to dominate, and these topics are highly correlated. Figure 5(a) illustrates the transition matrix inferred

by the PGDS, which is averaged over all time steps. Compared with the NS-PGDS, the PGDS can 298 not capture the informative time-varying transition dynamics. We also analyzed the features of the 299 proposed Dirichlet Markov chains. The left column of Figure 6 demonstrates transition matrices 300 of the first four sub-intervals of ICEWS dataset inferred by the NS-PGDS (Dir-Dir). Because of 301 the Dir-Dir construction, the consecutive transition matrices smoothly change over time and thus 302 the NS-PGDS may lack sufficient flexibility to capture rapid dynamics. The middle column of 303 Figure 6 illustrates the transition matrices inferred by the NS-PGDS (Dir-Gam-Dir), which takes 304 mutations among latent components into account and captured more complicated patterns. Transition 305 matrices inferred by the PR-Gam-Dir construction are shown in the right column of Figure 6, these 306 matrices not only exhibited sufficient flexibility but also captured sparser patterns compared with the 307 Dir-Gam-Dir construction.



Figure 6: From top to bottom are the first four transition matrices inferred by different Dirichlet Markov chains from ICEWS dataset. Top row: Matrices inferred by the Dir-Dir construction. Middle row: Matrices inferred by the Dir-Gam-Dir construction. Bottom row: Matrices inferred by the PR-Gam-Dir construction.

309 7 Conclusion

308

The Poisson-gamma dynamical systems with time-varying transition matrices, have been proposed to 310 capture complicated dynamics observed in *non-stationary* count sequences. In particular, Dirichlet 311 Markov chains are constructed to allow the underlying transition matrices to evolve over time. 312 Although the Dirichlet Markov processes lack conjugacy, we have developed tractable-but-efficient 313 Gibbs sampling algorithms to perform posterior simulation. The experiment results demonstrate the 314 improved performance of the proposed NS-PGDS in data smoothing and forecasting tasks, compared 315 with the PGDS with a stationary transition kernel. Moreover, the experimental results on several 316 real-world data sets show the explainable structures inferred by the proposed NS-PGDS. For the 317 future work, we plan to design a method that can find the point of change and thus the length of each 318 sub-interval can be determined automatically instead of a constant. We also consider to generalize 319 Dirichlet belief networks by incorporating the proposed Dirichlet Markov chain constructions, which 320 allow the hierarchical topics to mutate across layers, and thus can generate more rich text information. 321 And we also consider to capture non-stationary interaction dynamics among individuals over online 322 social networks in the future research. 323

324 **References**

- [1] David M Blei and John D Lafferty. Dynamic topic models. In *Proceedings of the 23rd international conference on Machine learning*, pages 113–120, 2006.
- [2] Xuerui Wang and Andrew McCallum. Topics over time: a non-markov continuous-time model of topical trends. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 424–433, 2006.
- [3] Patrick Jähnichen, Florian Wenzel, Marius Kloft, and Stephan Mandt. Scalable generalized
 dynamic topic models. In *International Conference on Artificial Intelligence and Statistics*,
 pages 1427–1435, 2018.
- [4] Daniel Sheldon and Thomas G Dietterich. Collective graphical models. In *Proceedings of the* 24th International Conference on Neural Information Processing Systems, pages 1161–1169,
 2011.
- [5] James Raymer, Arkadiusz Wiśniowski, Jonathan J Forster, Peter WF Smith, and Jakub Bijak.
 Integrated modeling of european migration. *Journal of the American Statistical Association*, 108(503):801–819, 2013.
- [6] Tom Wilson. Methods for estimating sub-state international migration: The case of australia.
 Spatial Demography, 5(3):171–192, 2017.
- [7] Philippe Wanner. How well can we estimate immigration trends using google data? *Quality & Quantity*, 55(4):1181–1202, 2021.
- [8] R. E. Kalman. A New Approach to Linear Filtering and Prediction Problems. *Journal of Basic Engineering*, 82(1):35–45, 1960.
- [9] Zoubin Ghahramani and Sam T Roweis. Learning nonlinear dynamical systems using an
 em algorithm. In *Proceedings of the 11th International Conference on Neural Information Processing Systems*, pages 431–437, 1998.
- [10] M Zhou and L Carin. Augment-and-conquer negative binomial processes. *Advances in Neural Information Processing Systems*, 4:2546–2554, 2012.
- [11] Mingyuan Zhou and Lawrence Carin. Negative binomial process count and mixture modeling.
 IEEE Transactions on Pattern Analysis & Machine Intelligence, 37(02):307–320, 2015.
- [12] Aaron Schein, John Paisley, David M Blei, and Hanna Wallach. Bayesian poisson tensor
 factorization for inferring multilateral relations from sparse dyadic event counts. In *Proceedings* of the 21th ACM SIGKDD International conference on knowledge discovery and data mining,
 pages 1045–1054, 2015.
- [13] Aaron Schein, Mingyuan Zhou, and Hanna Wallach. Poisson-gamma dynamical systems. In
 Proceedings of the 30th International Conference on Neural Information Processing Systems,
 pages 5012–5020, 2016.
- [14] Aaron Schein, Mingyuan Zhou, David Blei, and Hanna Wallach. Bayesian poisson tucker
 decomposition for learning the structure of international relations. In *International Conference on Machine Learning*, pages 2810–2819, 2016.
- [15] Ayan Acharya, Joydeep Ghosh, and Mingyuan Zhou. Nonparametric bayesian factor analysis
 for dynamic count matrices. In *Artificial Intelligence and Statistics*, pages 1–9, 2015.
- [16] Aaron Schein, Scott W Linderman, Mingyuan Zhou, David M Blei, and Hanna Wallach.
 Poisson-randomized gamma dynamical systems. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 782–793, 2019.
- [17] Jonathan Chang and David Blei. Relational topic models for document networks. In *Proceedings* of the Twelth International Conference on Artificial Intelligence and Statistics, pages 81–88,
 2009.

- [18] H Juliette T Unwin, Swapnil Mishra, Valerie C Bradley, Axel Gandy, Thomas A Mellan, Helen
 Coupland, Jonathan Ish-Horowicz, Michaela AC Vollmer, Charles Whittaker, Sarah L Filippi,
 et al. State-level tracking of covid-19 in the united states. *Nature Communications*, 11(1):1–9,
 2020.
- [19] Rainer Winkelmann. *Econometric Analysis of Count Data*. Springer Publishing Company,
 Incorporated, 5th edition, 2008.
- [20] Guy Grossman, Soojong Kim, Jonah M Rexer, and Harsha Thirumurthy. Political partisanship
 influences behavioral responses to governors' recommendations for covid-19 prevention in the
 united states. *Proceedings of the National Academy of Sciences*, 117(39):24144–24153, 2020.
- IHME COVID-19 Forecasting Team. Modeling covid-19 scenarios for the united states. *Nature medicine*, 27(1):94–105, 2021.
- [22] Luzhao Feng, Ting Zhang, Qing Wang, Yiran Xie, Zhibin Peng, Jiandong Zheng, Ying Qin,
 Muli Zhang, Shengjie Lai, Dayan Wang, et al. Impact of covid-19 outbreaks and interventions
 on influenza in china and the united states. *Nature communications*, 12(1):3249, 2021.
- [23] Shaobo Han, Lin Du, Esther Salazar, and Lawrence Carin. Dynamic rank factor model for text
 streams. In *Proceedings of the 27th International Conference on Neural Information Processing Systems-Volume 2*, pages 2663–2671, 2014.
- Rahi Kalantari and Mingyuan Zhou. Graph gamma process generalized linear dynamical
 systems. *arXiv preprint arXiv:2007.12852*, 2020.
- [25] Lin Yuan and John D Kalbfleisch. On the bessel distribution and related problems. *Annals of the Institute of Statistical Mathematics*, 52:438–447, 2000.
- ³⁹¹ [26] Yee Whye Teh, Michael I Jordan, Matthew J Beal, and David M Blei. Hierarchical dirichlet ³⁹² processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [27] Norman L Johnson, Adrienne W Kemp, and Samuel Kotz. *Univariate discrete distributions*,
 volume 444. John Wiley & Sons, 2005.
- [28] Mingyuan Zhou. Nonparametric bayesian negative binomial factor analysis. *Bayesian Analysis*, 13(4):1065–1093, 2018.
- [29] Chengyue Gong and Win-bin Huang. Deep dynamic poisson factorization model. In *Proceedings* of the 31st International Conference on Neural Information Processing Systems, pages 1665–
 1673, 2017.
- [30] Sikun Yang and Heinz Koeppl. Dependent relational gamma process models for longitudinal
 networks. In *International Conference on Machine Learning*, pages 5551–5560, 2018.
- [31] Dandan Guo, Bo Chen, Hao Zhang, and Mingyuan Zhou. Deep poisson gamma dynamical
 systems. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 8451–8461, 2018.
- [32] Wenchao Chen, Bo Chen, Yicheng Liu, Qianru Zhao, and Mingyuan Zhou. Switching poisson
 gamma dynamical systems. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 2029–2036, 2021.
- [33] Louis Filstroff, Olivier Gouvert, Cédric Févotte, and Olivier Cappé. A comparative study of
 gamma markov chains for temporal non-negative matrix factorization. *IEEE Transactions on Signal Processing*, 69:1614–1626, 2021.
- [34] Peter AW Lewis, Edward McKenzie, and David Kennedy Hugus. Gamma processes. *Stochastic Models*, 5(1):1–30, 1989.
- 413 [35] John Frank Charles Kingman. *Poisson processes*, volume 3. Clarendon Press, 1992.

414 A MCMC Inference

Notation. When expressing the full conditionals for Gibbs sampling, we use the shorthand "–" to denote all other variables. We use "•" as an index summation shorthand, e.g., $x_{\cdot j} = \sum_{i} x_{ij}$.

⁴¹⁷ In this section, we present a fully-conjugate and efficient Gibbs sampler for the proposed NS-PGDS.

The sampling algorithms depend on several key technical results, which we will repeatedly exploit, thus we list them below.

Negative-binomial Distribution. Let $y \sim \text{Pois}(c\lambda)$, and $\lambda \sim \text{Gam}(a, b)$. If we marginalize

421 over λ , then $y \sim \text{NB}\left(a, \frac{c}{b+c}\right)$ is a negative-binomial distributed random variable. We can further

parameterize it as $y \sim NB(a, g(\zeta))$, where $g(z) = 1 - \exp(-z)$ and $\zeta = \ln\left(1 + \frac{c}{b}\right)$.

Lemma 1. If $y \sim \text{NB}(a, g(\zeta))$ and $l \sim \text{CRT}(y, a)$, where $\text{CRT}(\cdot)$ represents Chinese restaurant table distribution [26], then the joint distribution of y and l can be equivalently distributed as $y \sim \text{SumLog}(l, g(\zeta))$ and $l \sim \text{Pois}(a\zeta)$ [11], i.e.

 $NB(y; a, g(\zeta)) CRT(l; y, a) = SumLog(y; l, g(\zeta)) Pois(l; a\zeta),$

where SumLog $(l, g(\zeta)) = \sum_{i=1}^{l} x_i$ and $x_i \sim \text{Log}(g(\zeta))$ are independently and identically logarithmic distributed random variables [27].

Lemma 2. Suppose $\mathbf{n} = (n_1, \dots, n_K)$ and $\mathbf{n} \mid n \sim \text{DirMult}(n, r_1, \dots, r_K)$, where $\text{DirMult}(\cdot)$ refers to Dirichlet-multimonial distribution. We sample the augmented variable $q \mid n \sim \text{Beta}(n, r_{\cdot})$, where $r_{\cdot} = \sum_{k=1}^{K} r_k$. According to [28], conditioning on q, we have $n_k \sim \text{NB}(r_k, q)$.

Lemma 3. If $y_{\cdot} = \sum_{s=1}^{S} y_s$, and $y_s \stackrel{\text{i.i.d}}{\sim} \operatorname{Pois}(\lambda_s), s = 1, \dots, S$. Then $y_{\cdot} \sim \operatorname{Pois}(\sum_{s=1}^{S} \lambda_s)$ and $(y_1, \dots, y_S) \sim \operatorname{Mult}(y_{\cdot}, (\frac{\lambda_1}{\sum_{s=1}^{S} \lambda_s}, \dots, \frac{\lambda_S}{\sum_{s=1}^{S} \lambda_s}))$, where Mult (\cdot) represents multinomial distribution [35].

Sampling $y_{vk}^{(t)}$: Use the relationship between Poisson and multinomial distributions as described by Lemma 3, given observed counts and latent parameters, we sample

$$\left(\left(y_{vk}^{(t)}\right)_{k=1}^{K} \mid -\right) \sim \operatorname{Mult}\left(y_{v}^{(t)}, \left(\frac{\phi_{vk}\theta_{k}^{(t)}}{\sum_{k=1}^{K}\phi_{vk}\theta_{k}^{(t)}}\right)_{k=1}^{K}\right).$$
(6)

Then the distribution of $y_{vk}^{(t)}$ is $y_{vk}^{(t)} \sim \text{Pois}(\delta^{(t)}\phi_{vk}\theta_k^{(t)})$.

434 Sampling ϕ_k : Via Dirichlet-multinomial conjugacy, the posterior of ϕ_k is

$$(\boldsymbol{\phi}_{\boldsymbol{k}} \mid -) \sim \operatorname{Dir}\left(\boldsymbol{\epsilon}_{0} + \sum_{t=1}^{T} y_{1k}^{(t)}, \cdots, \boldsymbol{\epsilon}_{0} + \sum_{t=1}^{T} y_{Vk}^{(t)}\right).$$
(7)

435 **Marginalizing over** $\theta_k^{(t)}$: Note that $y_v^{(t)} = y_{v}^{(t)} = \sum_{k=1}^K y_{vk}^{(t)}$ and $y_{vk}^{(t)} \sim \operatorname{Pois}(\delta^{(t)}\phi_{vk}\theta_k^{(t)})$. Then 436 we define $y_{\cdot k}^{(t)} = \sum_{v=1}^V y_{vk}^{(t)}$. Because $\sum_{v=1}^V \phi_{vk} = 1$, we obtain $y_{\cdot k}^{(t)} \sim \operatorname{Pois}(\delta^{(t)}\theta_k^{(t)})$.

437 We start by marginalizing over $\theta_k^{(T)}$, using the definition of negative-binomial distribution, we obtain

$$y_{\cdot k}^{(T)} \sim \text{NB}\left(\tau_0 \sum_{k_2=1}^{K} \pi_{kk_2}^{i(T-1)} \theta_{k_2}^{(T-1)}, g\left(\zeta^{(T)}\right)\right),$$

where $\zeta^{(T)} = \ln(1 + \frac{\delta^{(T)}}{\tau_0})$. Next, we further marginalize over $\theta_k^{(T-1)}$. To this end, we first sample auxiliary variables

$$l_{k}^{(T)} \sim \operatorname{CRT}\left(y_{\cdot k}^{(T)}, \tau_{0} \sum_{k_{2}=1}^{K} \pi_{k k_{2}}^{i(T-1)} \theta_{k_{2}}^{(T-1)}\right)$$

⁴⁴⁰ By Lemma 1, the joint distribution of $y_{\cdot k}^{(T)}$ and $l_k^{(T)}$ can be expressed as

$$y_{k}^{(T)} \sim \text{SumLog}\left(l_{k}^{(T)}, g\left(\zeta^{(T)}\right)\right) \text{ and } l_{k}^{(T)} \sim \text{Pois}\left(\zeta^{(T)}\tau_{0}\sum_{k_{2}=1}^{K} \pi_{kk_{2}}^{i(T-1)}\theta_{k_{2}}^{(T-1)}\right).$$

441 Via Lemma 3, we re-express the auxiliary variables as

$$l_{k}^{(T)} = l_{k}^{(T)} = \sum_{k_{2}=1}^{K} l_{kk_{2}}^{(T)}, \text{ and obtain } l_{kk_{2}}^{(T)} \sim \operatorname{Pois}\left(\zeta^{(T)}\tau_{0}\pi_{kk_{2}}^{i(T-1)}\theta_{k_{2}}^{(T-1)}\right).$$

Then we define $l_{\cdot k}^{(T)} = \sum_{k_1=1}^{K} l_{k_1 k}^{(T)}$. Leveraging Lemma 3 and $\sum_{k_1=1}^{K} \pi_{k_1 k}^{i(T-1)} = 1$, we obtain $l_{\cdot k}^{(T)} \sim \operatorname{Pois}\left(\zeta^{(T)} \tau_0 \theta_k^{(T-1)}\right)$ and $\left(l_{1k}^{(T)}, \cdots, l_{Kk}^{(T)}\right) \sim \operatorname{Mult}\left(l_{\cdot k}^{(T)}, \left(\pi_{1k}^{i(T-1)}, \cdots, \pi_{Kk}^{i(T-1)}\right)\right)$.

⁴⁴³ Next, note that $y_{\cdot k}^{(T-1)} \sim \text{Pois}(\delta^{(T-1)}\theta_k^{(T-1)})$, if we introduce $m_k^{(T-1)} = y_{\cdot k}^{(T-1)} + l_{\cdot k}^{(T)}$, then we ⁴⁴⁴ have

$$m_k^{(T-1)} \sim \operatorname{Pois}\left(\theta_k^{(T-1)}\left(\delta^{(T-1)} + \zeta^{(T)}\tau_0\right)\right).$$

Because the prior of $\theta_k^{(T-1)}$ is gamma distributed, by the definition of negative-binomial distribution, we can again marginalize over $\theta_k^{(T-1)}$ to obtain

$$m_k^{(T-1)} \sim \text{NB}\left(\tau_0 \sum_{k_2=1}^K \pi_{kk_2}^{i(T-2)} \theta_{k_2}^{(T-2)}, g\left(\zeta^{(T-1)}\right)\right)$$

447 where $\zeta^{(T-1)} = \ln(1 + \frac{\delta^{(T-1)}}{\tau_0} + \zeta^{(T)})$. Then we introduce auxiliary variables

$$l_k^{(T-1)} \sim \text{CRT}\left(m_k^{(T-1)}, \tau_0 \sum_{k_2=1}^K \pi_{kk_2}^{i(T-2)} \theta_{k_2}^{(T-2)}\right).$$

And similar to the case for t = T, we can obtain

$$l_{\cdot k}^{(T-1)} \sim \operatorname{Pois}\left(\zeta^{(T-1)} \tau_0 \theta_k^{(T-2)}\right) \text{ and } m_k^{(T-2)} \sim \operatorname{NB}\left(\tau_0 \sum_{k_2=1}^K \pi_{kk_2}^{i(T-3)} \theta_{k_2}^{(T-3)}, g\left(\zeta^{(T-2)}\right)\right).$$

Thus we have marginalized over $\theta_k^{(T-2)}$. Note that we can repeat this marginalization process recursively until t = 1 with $\zeta^{(t)} = \ln(1 + \frac{\delta^{(t)}}{\tau_0} + \zeta^{(t+1)})$ and $m_k^{(T)} = y_{\cdot k}^{(T)}$ to maginalize over all the $\theta_k^{(t)}$.

Sampling $\theta_k^{(t)}$: Via the above marginalization process, to sample from the posterior of $\theta_k^{(t)}$, we first sample the auxiliary variables. Setting $l_{\cdot k}^{(T+1)} = 0$ and $\zeta^{(T+1)} = 0$, we sample the augmented variables backwards from $t = T, \dots, 2$,

$$\left(l_{k}^{(t)} \mid -\right) \sim \operatorname{CRT}\left(y_{\cdot k}^{(t)} + l_{\cdot k}^{(t+1)}, \tau_0 \sum_{k_2=1}^{K} \pi_{k k_2}^{i(t-1)} \theta_{k_2}^{(t-1)}\right),\tag{8}$$

$$\left(l_{k1}^{(t)}, \cdots, l_{kK}^{(t)} \mid - \right) \sim \operatorname{Mult} \left(l_{k\cdot}^{(t)}, \left(\frac{\pi_{k1}^{i(t-1)} \theta_1^{(t-1)}}{\sum_{k_2=1}^K \pi_{kk_2}^{i(t-1)} \theta_{k_2}^{(t-1)}}, \cdots, \frac{\pi_{kK}^{i(t-1)} \theta_K^{(t-1)}}{\sum_{k_2=1}^K \pi_{kk_2}^{i(t-1)} \theta_{k_2}^{(t-1)}} \right) \right).$$
(9)

455 And via Lemma 3, we obtain

$$\left(l_{1k}^{(t)}, \cdots, l_{Kk}^{(t)}\right) \sim \text{Mult}\left(l_{\cdot k}^{(t)}, \pi_{1k}^{i(t-1)}, \cdots, \pi_{Kk}^{i(t-1)}\right)$$
 (10)

456 We compute $\zeta^{(t)}$ recursively via

$$\zeta^{(t)} = \ln\left(1 + \frac{\delta^{(t)}}{\tau_0} + \zeta^{(t+1)}\right).$$
(11)

After sampling the auxiliary variables, then for $t = 1, \dots, T$, by Poisson-gamma conjugacy, we obtain

$$\left(\theta_{k}^{(1)} \mid -\right) \sim \operatorname{Gam}\left(y_{.k}^{(1)} + l_{.k}^{(2)} + \tau_{0}\nu_{k}, \tau_{0} + \delta^{(1)} + \zeta^{(2)}\tau_{0}\right),\tag{12}$$

$$\left(\theta_{k}^{(t)} \mid -\right) \sim \operatorname{Gam}\left(y_{\cdot k}^{(t)} + l_{\cdot k}^{(t+1)} + \tau_{0} \sum_{k_{2}=1}^{K} \pi_{k k_{2}}^{i(t-1)} \theta_{k_{2}}^{(t-1)}, \tau_{0} + \delta^{(t)} + \zeta^{(t+1)} \tau_{0}\right).$$
(13)

459

Sampling $\Pi^{(i)}$: We define M as the length of each sub-interval, and I as the number of intervals. For i = I, by Eq.(10), $(l_{1k}^{(I)}, \dots, l_{Kk}^{(I)})$ is multinomial distributed. Thus by multinomial-Dirichlet conjugacy, we obtain

$$\left(\pi_{k}^{(I)}\mid-\right)\sim\operatorname{Dir}\left(\alpha_{1k}^{(I)}+l_{1k}^{(I)},\cdots,\alpha_{Kk}^{(I)}+l_{Kk}^{(I)}\right),$$
(14)

where $l_{k_1k}^{(I)}$ indicates the summation of $l_{k_1k}^{(t)}$ over *I*-th sub-interval, i.e. $l_{k_1k}^{(I)} = \sum_{t=(I-1)M+1}^{T} l_{k_1k}^{(t)}$

Inference for Dirichlet-Dirichlet Markov chains. For Dirichlet-Dirichlet Markov chains, $\alpha_{k_1k}^{(i)} = \eta K \pi_{k_1k}^{(i-1)}$. By Eq.(10), $(l_{1k}^{(i)}, \dots, l_{Kk}^{(i)})$ is multinomial distributed. If we marginalize $(\pi_{1k}^{(i)}, \dots, \pi_{Kk}^{(i)})$, ($l_{1k}^{(i)}, \dots, l_{Kk}^{(i)}$) will be Dirichlet-multinomial distributed. Thus by Lemma 2, for i = I, we first sample the auxiliary variables as

$$\left(q_{k}^{(I)}\mid-\right)\sim\operatorname{Beta}\left(l_{\cdot k}^{(I)},\eta K\right) \text{ and } \left(h_{k_{1}k}^{(I)}\mid-\right)\sim\operatorname{CRT}\left(l_{k_{1}k}^{(I)},\eta K\pi_{k_{1}k}^{(I-1)}\right).$$
(15)

Similarly, by Eq.(18), $(h_{1k}^{(i)}, \dots, h_{Kk}^{(i)})$ is also Dirichlet-multinomial distributed. Thus for $i = I - 1, \dots, 2$, we sample the auxiliary variables as

$$\left(q_{k}^{(i)}\mid-\right)\sim\operatorname{Beta}\left(l_{\cdot k}^{(i)}+h_{\cdot k}^{(i+1)},\eta K\right) \text{ and } \left(h_{k_{1}k}^{(i)}\mid-\right)\sim\operatorname{CRT}\left(l_{k_{1}k}^{(i)}+h_{k_{1}k}^{(i+1)},\eta K\pi_{k_{1}k}^{(i-1)}\right),$$
(16)

where $l_{k_1k}^{(i)} = \sum_{i=1}^{iM} l_{k_1k}^{(t)}$ refers to the summation of $l_{k_1k}^{(t)}$ over *i*-th interval. Via Lemma 2, conditioning on $q_k^{(i)}$, we have

$$\left(l_{k_1k}^{(i)} + h_{k_1k}^{(i+1)}\right) \sim \text{NB}\left(\eta K \pi_{k_1k}^{(i-1)}, q_k^{(i)}\right)$$

472 Then via Lemma 1, we obtain

$$h_{k_1k}^{(i)} \sim \text{Pois}\left(-\eta K \pi_{k_1k}^{(i-1)} \ln\left(1 - q_k^{(i)}\right)\right).$$
 (17)

⁴⁷³ Note that by Eq.(17), $h_{k_1k}^{(i)}$ is Poisson distributed and by Lemma 3, we obtain

$$\left(h_{1k}^{(i)}, \cdots, h_{Kk}^{(i)}\right) \sim \operatorname{Mult}\left(h_{\cdot k}^{(i)}, \left(\pi_{1k}^{(i-1)}, \cdots, \pi_{Kk}^{(i-1)}\right)\right).$$
 (18)

474 In addition, note that

$$\left(l_{1k}^{(i-1)}, \cdots, l_{Kk}^{(i-1)}\right) \sim \operatorname{Mult}\left(l_{\cdot k}^{(i-1)}, \left(\pi_{1k}^{(i-1)}, \cdots, \pi_{Kk}^{(i-1)}\right)\right),$$

via Dirichlet-multinomial conjugacy, for $i = I - 1, \cdots, 2$, we obtain

$$\left(\boldsymbol{\pi}_{k}^{(i)} \mid -\right) \sim \operatorname{Dir}\left(\eta K \pi_{1k}^{(i-1)} + l_{1k}^{(i)} + h_{1k}^{(i+1)}, \cdots, \eta K \pi_{Kk}^{(i-1)} + l_{Kk}^{(i)} + h_{Kk}^{(i+1)}\right).$$
(19)

476 Specifically, for i = 1, we have

$$\left(\boldsymbol{\pi}_{k}^{(1)} \mid -\right) \sim \operatorname{Dir}\left(\nu_{1}\nu_{k} + l_{1k}^{(1)} + h_{1k}^{(2)}, \cdots, \xi\nu_{k} + l_{kk}^{(1)} + h_{kk}^{(2)}, \cdots, \nu_{K}\nu_{k} + l_{Kk}^{(1)} + h_{Kk}^{(2)}\right).$$
(20)

For sampling η , note that $(h_{k_1k}^{(i)} \mid -) \sim \text{Pois}(-\eta K \pi_{k_1k}^{(i-1)} \ln\left(1 - q_k^{(i)}\right))$, $i = I, \dots, 2$. Given the prior $\eta \sim \text{Gam}(e_0, f_0)$, via Poisson-gamma conjugacy, we obtain

$$(\eta \mid -) \sim \operatorname{Gam}\left(e_0 + \sum_{i=2}^{I} \sum_{k_1=1}^{K} \sum_{k_2=1}^{K} h_{k_1 k_2}^{(i)}, f_0 - K \sum_{i=2}^{I} \sum_{k=1}^{K} \ln\left(1 - q_k^{(i)}\right)\right).$$
(21)

Inference for Dirichlet-Gamma-Dirichlet Markov chains. For Dirichlet-Gamma-Dirichlet Markov
 chains

$$\alpha_{k_1k}^{(i)} \sim \operatorname{Gam}\left(\gamma_k^{(i-1)} \sum_{k2=1}^K \psi_{kk_1k_2}^{(i-1)} \pi_{k_2k}^{(i-1)}, c_k^{(i)}\right).$$

By Eq.(10), $(l_{1k}^{(i)}, \dots, l_{Kk}^{(i)})$ is multinomial distributed. If we marginalize $(\pi_{1k}^{(i)}, \dots, \pi_{Kk}^{(i)})$, $(l_{1k}^{(i)}, \dots, l_{Kk}^{(i)})$ will be Dirichlet-multinomial distributed. Thus by Lemma 2, for i = I, we first sample the auxiliary variables as 481 482 483

$$\left(q_{k}^{(I)}\mid-\right)\sim\operatorname{Beta}\left(l_{\cdot k}^{(I)},\alpha_{\cdot k}^{(I)}\right) \text{ and } \left(h_{k_{1}k}^{(I)}\mid-\right)\sim\operatorname{CRT}\left(l_{k_{1}k}^{(I)},\alpha_{k_{1}k}^{(I)}\right).$$
(22)

Similarly, by Eq.(27), $(g_{\cdot 1k}^{(i)}, \cdots, g_{\cdot Kk}^{(i)})$ is also Dirichlet-multinomial distributed. Thus for $i = I - 1, \cdots, 2$, we sample the auxiliary variables as 484 485

$$\left(q_{k}^{(i)} \mid -\right) \sim \text{Beta}\left(l_{\cdot k}^{(i)} + g_{\cdot k}^{(i+1)}, \alpha_{\cdot k}^{(i)}\right) \text{ and } \left(h_{k_{1}k}^{(i)} \mid -\right) \sim \text{CRT}\left(l_{k_{1}k}^{(i)} + g_{\cdot k_{1}k}^{(i+1)}, \alpha_{k_{1}k}^{(i)}\right).$$
(23)

Via Lemma 2, conditioning on $q_k^{(i)}$, we have 486

$$\left(l_{k_1k}^{(i)} + g_{\cdot k_1k}^{(i+1)}\right) \sim \text{NB}\left(\alpha_{k_1k}^{(i)}, q_k^{(i)}\right)$$

Then via Lemma 1, we obtain 487

$$h_{k_1k}^{(i)} \sim \operatorname{Pois}\left(-\alpha_{k_1k}^{(i)}\ln\left(1-q_k^{(i)}\right)\right)$$

Thus via Poisson-gamma conjugacy, we obtain 488

$$\left(\alpha_{k_{1}k}^{(i)}\mid-\right)\sim\operatorname{Gam}\left(\gamma_{k}^{(i-1)}\sum_{k_{2}=1}^{K}\psi_{kk_{1}k_{2}}^{(i-1)}\pi_{k_{2}k}^{(i-1)}+h_{k_{1}k}^{(i)},c_{k}^{(i)}-\ln\left(1-q_{k}^{(i)}\right)\right).$$
(24)

Marginalizing over $\alpha_{k_1k}^{(i)}$, and via the definition of negative-binomial distribution, we have 489

$$h_{k_1k}^{(i)} \sim \text{NB}\left(\gamma_k^{(i-1)} \sum_{k2=1}^K \psi_{kk_1k_2}^{(i-1)} \pi_{k_2k}^{(i-1)}, \frac{-\ln\left(1-q_k^{(i)}\right)}{c_k^{(i)} - \ln\left(1-q_k^{(i)}\right)}\right)$$

Then using Lemma 1, we sample 490

$$\left(g_{k_1k}^{(i)} \mid -\right) \sim \operatorname{CRT}\left(h_{k_1k}^{(i)}, \gamma_k^{(i-1)} \sum_{k_{2}=1}^{K} \psi_{kk_1k_2}^{(i-1)} \pi_{k_2k}^{(i-1)}\right),\tag{25}$$

and obtain 491

$$g_{k_1k}^{(i)} \sim \operatorname{Pois}\left(\gamma_k^{(i-1)} \sum_{k_{2=1}}^K \psi_{kk_1k_2}^{(i-1)} \pi_{k_2k}^{(i-1)} \ln\left(1 - \ln\left(1 - q_k^{(i)}\right) / c_k^{(i)}\right)\right)$$

If we define $g_{k_1k}^{(i)} = g_{k_1 \cdot k}^{(i)} = \sum_{k2=1}^{K} g_{k_1k_2k}^{(i)}$, and augment 492

$$\left(g_{k_{1}1k}^{(i)}, \cdots, g_{k_{1}Kk}^{(i)}\right) \sim \operatorname{Mult}\left(g_{k_{1}k}^{(i)}, \left(\psi_{kk_{1}k_{2}}^{(i-1)} \pi_{k_{2}k}^{(i-1)}\right)_{k_{2}=1}^{K}\right).$$
(26)

By Lemma 3, we have 493

$$g_{k_1k_2k}^{(i)} \sim \operatorname{Pois}\left(\gamma^{(i-1)}\psi_{kk_1k_2}^{(i-1)}\pi_{k_2k}^{(i-1)}\ln\left(1-\ln\left(1-q_k^{(i)}\right)/c_k^{(i)}\right)\right)$$

Using Lemma 3 and $\sum_{k_1}^{K} \psi_{kk_1k_2}^{(i-1)} = 1$, we have, 494

$$\begin{pmatrix} g_{\cdot 1k}^{(i)}, \cdots, g_{\cdot Kk}^{(i)} \end{pmatrix} \sim \operatorname{Mult} \left(g_{\cdot k}^{(i)}, \left(\pi_{k_1 k}^{(i-1)} \right)_{k_1 = 1}^K \right),$$

$$g_{1k_2 k}^{(i)}, \cdots, g_{Kk_2 k}^{(i)} \right) \sim \operatorname{Mult} \left(g_{\cdot k_2 k}^{(i)}, \left(\psi_{kk_1 k_2}^{(i-1)} \right)_{k_1 = 1}^K \right).$$

$$(27)$$

495

$$\left(g_{1k_{2k}}^{(i)},\cdots,g_{Kk_{2k}}^{(i)}\right) \sim \operatorname{Mult}\left(g_{\cdot k_{2k}}^{(i)},\left(\psi_{kk_{1}k_{2}}^{(i-1)}\right)_{k_{1}=1}^{K}\right).$$

Thus by Dirichlet-multinomial conjugacy, for $i = I, \dots, 2$, we can obtain 496

$$\begin{pmatrix} \left(\psi_{k1k_{2}}^{(i-1)}, \cdots, \psi_{kKk_{2}}^{(i-1)}\right) \mid - \end{pmatrix} \sim \operatorname{Dir} \left(\epsilon_{0} + g_{1k_{2}k}^{(i)}, \cdots, \epsilon_{0} + g_{Kk_{2}k}^{(i)}\right),$$

$$\begin{pmatrix} \pi_{k}^{(i-1)} \mid - \end{pmatrix} \sim \operatorname{Dir} \left(\alpha_{1k}^{(i-1)} + l_{1k}^{(i-1)} + g_{\cdot 1k}^{(i)}, \cdots, \alpha_{Kk}^{(i-1)} + l_{Kk}^{(i-1)} + g_{\cdot Kk}^{(i)}\right).$$

$$(28)$$

497 For sampling $\gamma_k^{(i-1)}$, note that by Eq.(26) and $\sum_{k_1}^K \psi_{kk_1k_2}^{(i-1)} = 1$, we have

$$g_{\cdot k}^{(i)} = \sum_{k_1=1}^{K} g_{k_1 k}^{(i)} \text{ and } g_{\cdot k}^{(i)} \sim \operatorname{Pois}\left(\gamma_k^{(i-1)} \ln\left(1 - \ln\left(1 - q_k^{(i)}\right) / c_k^{(i)}\right)\right).$$
(30)

⁴⁹⁸ Thus via Poisson-gamma conjugacy, we obtain

$$\left(\gamma_k^{(i-1)} \mid -\right) \sim \operatorname{Gam}\left(\epsilon_0 + g_{\cdot k}^{(i)}, \epsilon_0 + \ln\left(1 - \ln\left(1 - q_k^{(i)}\right)\right)\right). \tag{31}$$

⁴⁹⁹ By gamma-gamma conjugacy, we have

$$\left(c_k^{(i)} \mid -\right) \sim \operatorname{Gam}\left(\epsilon_0 + \gamma_k^{(i-1)}, \epsilon_0 + \sum_{k_1=1}^K \alpha_{k_1k}^{(i)}\right).$$
(32)

Inference for Dirichlet-Randomized-Gamma-Dirichlet Markov chains. For Dirichlet Randomized-Gamma-Dirichlet Markov chains,

$$\alpha_{k_1k}^{(i)} \sim \mathrm{RG1}\left(\epsilon^{\alpha}, \gamma^{(i-1)} \sum_{k_2=1}^{K} \psi_{kk_1k_2}^{(i-1)} \pi_{k_2k}^{(i-1)}, c_k^{(i)}\right),$$

⁵⁰² which can be equivalently represented as

$$\alpha_{k_1k}^{(i)} \sim \operatorname{Gam}\left(g_{k_1k}^{(i)} + \epsilon^{\alpha}, c_k^{(i)}\right), \text{ and } g_{k_1k}^{(i)} = \operatorname{Pois}\left(\gamma^{(i-1)} \sum_{k_2=1}^{K} \psi_{kk_1k_2}^{(i-1)} \pi_{k_2k}^{(i-1)}\right).$$

⁵⁰³ By Eq.(10), $(l_{1k}^{(i)}, \dots, l_{Kk}^{(i)})$ is multinomial distributed. If we marginalize $(\pi_{1k}^{(i)}, \dots, \pi_{Kk}^{(i)})$, ⁵⁰⁴ $(l_{1k}^{(i)}, \dots, l_{Kk}^{(i)})$ will be Dirichlet-multinomial distributed. Thus by Lemma 2, for i = I, we first ⁵⁰⁵ sample the auxiliary variables as

$$\left(q_{k}^{(I)}\mid-\right)\sim\operatorname{Beta}\left(l_{\cdot k}^{(I)},\alpha_{\cdot k}^{(I)}\right) \text{ and } \left(h_{k_{1}k}^{(I)}\mid-\right)\sim\operatorname{CRT}\left(l_{k_{1}k}^{(I)},\alpha_{k_{1}k}^{(I)}\right).$$
(33)

Similarly, by Eq.(39), $(g_{.1k}^{(i)}, \dots, g_{.Kk}^{(i)})$ is also Dirichlet-multinomial distributed. Thus for $i = I - 1, \dots, 2$, we sample the auxiliary variables as

$$\left(q_{k}^{(i)}\mid-\right)\sim\operatorname{Beta}\left(l_{\cdot k}^{(i)}+g_{\cdot k}^{(i+1)},\alpha_{\cdot k}^{(i)}\right) \text{ and } \left(h_{k_{1}k}^{(i)}\mid-\right)\sim\operatorname{CRT}\left(l_{k_{1}k}^{(i)}+g_{\cdot k_{1}k}^{(i+1)},\alpha_{k_{1}k}^{(i)}\right).$$
(34)

508 Via Lemma 2, conditioning on $q_k^{(i)}$, we have

$$\left(l_{k_1k}^{(i)} + g_{\cdot k_1k}^{(i+1)}\right) \sim \operatorname{NB}\left(\alpha_{k_1k}^{(i)}, q_k^{(i)}\right)$$

509 Then via Lemma 1, we obtain

$$h_{k_1k}^{(i)} \sim \operatorname{Pois}\left(-\alpha_{k_1k}^{(i)} \ln\left(1 - q_k^{(i)}\right)\right)$$

510 Via Poisson-gamma conjugacy, we first sample

$$\left(\alpha_{k_{1}k}^{(i)}\mid-\right)\sim\operatorname{Gam}\left(g_{k_{1}k}^{(i)}+\epsilon^{\alpha}+h_{k_{1}k}^{(i)},c_{k}^{(i)}-\ln\left(1-q_{k}^{(i)}\right)\right).$$
(35)

511 If $\epsilon^{\alpha} > 0$, we can sample the posterior of $g_{k_1k}^{(i)}$ via

$$\left(g_{k_{1}k}^{(i)}\mid-\right) \sim \text{Bessel}\left(\epsilon^{\alpha}-1, 2\sqrt{\alpha_{k_{1}k}^{(i)}c_{k}^{(i)}\gamma_{k}^{(i-1)}\sum_{k=1}^{K}\psi_{kk_{1}k_{2}}^{(i-1)}\pi_{k_{2}k}^{(i-1)}}\right),\tag{36}$$

where $\operatorname{Bessel}(\cdot)$ denotes Bessel distribution. If $\epsilon^{\alpha} = 0$, we sample $g_{k_1k}^{(i)}$ via

$$\begin{pmatrix} g_{k_1k}^{(i)} \mid - \end{pmatrix} \sim \begin{cases} \operatorname{Pois}\left(\frac{c_k^{(i)}\gamma_k^{(i-1)}\sum_{k=1}^{K}\psi_{kk_1k_2}^{(i-1)}\pi_{k_2k}^{(i-1)}}{c_k^{(i)} - \ln(1-q_k^{(i)})}\right) & \text{if } h_{k_1k}^{(i)} = 0\\ \operatorname{SCH}\left(h_{k_1k}^{(i)}, \frac{c_k^{(i)}\gamma_k^{(i-1)}\sum_{k=1}^{K}\psi_{kk_1k_2}^{(i-1)}\pi_{k_2k}^{(i-1)}}{c_k^{(i)} - \ln(1-q_k^{(i)})}\right) & \text{otherwise,} \end{cases}$$
(37)

where SCH (\cdot) denotes the shifted confluent hypergeometric distribution [16].

514 Defining $g_{k_1k}^{(i)} = g_{k_1 \cdot k}^{(i)} = \sum_{k_2=1}^{K} g_{k_1k_2k}^{(i)}$, we first augment

$$\left(g_{k_11k}^{(i)}, \cdots, g_{k_1Kk}^{(i)}\right) \sim \operatorname{Mult}\left(g_{k_1k}^{(i)}, \left(\psi_{kk_1k_2}^{(i-1)} \pi_{k_2k}^{(i-1)}\right)_{k_2=1}^K\right).$$
(38)

515 By Lemma 3, we have

$$g_{k_1k_2k}^{(i)} \sim \operatorname{Pois}\left(\gamma^{(i-1)}\psi_{kk_1k_2}^{(i-1)}\pi_{k_2k}^{(i-1)}\right),$$

and because $\sum_{k_1}^{K} \psi_{kk_1k_2}^{(i-1)} = 1$, we have

$$\left(g_{\cdot 1k}^{(i)}, \cdots, g_{\cdot Kk}^{(i)}\right) \sim \operatorname{Mult}\left(g_{\cdot k}^{(i)}, \left(\pi_{k_1 k}^{(i-1)}\right)_{k_1 = 1}^K\right), \text{ and}$$

$$\left(g_{1 k_2 k}^{(i)}, \cdots, g_{K k_2 k}^{(i)}\right) \sim \operatorname{Mult}\left(g_{\cdot k_2 k}^{(i)}, \left(\psi_{k k_1 k_2}^{(i-1)}\right)_{k_1 = 1}^K\right).$$

$$(39)$$

Thus by Dirichlet-multinomial conjugacy, for $i = I, \dots, 2$, we have

$$\left(\left(\psi_{k1k_{2}}^{(i-1)}, \cdots, \psi_{kKk_{2}}^{(i-1)}\right) \mid -\right) \sim \operatorname{Dir}\left(\epsilon_{0} + g_{1k_{2}k}^{(i)}, \cdots, \epsilon_{0} + g_{Kk_{2}k}^{(i)}\right),\tag{40}$$

518

$$\left(\pi_{k}^{(i-1)} \mid -\right) \sim \operatorname{Dir}\left(\alpha_{1k}^{(i-1)} + l_{1k}^{(i-1)} + g_{\cdot 1k}^{(i)}, \cdots, \alpha_{Kk}^{(i-1)} + l_{Kk}^{(i-1)} + g_{\cdot Kk}^{(i)}\right).$$
(41)

519 Via Poisson-gamma conjugacy, we obtain

$$\left(\gamma_k^{(i-1)} \mid -\right) \sim \operatorname{Gam}\left(\epsilon_0 + g_{\cdot k}^{(i)}, \epsilon_0 + 1\right).$$
(42)

520 By gamma-gamma conjugacy, we have

$$\left(c_{k}^{(i)}\mid-\right)\sim\operatorname{Gam}\left(\epsilon_{0}+\gamma_{k}^{(i-1)},\epsilon_{0}+\sum_{k_{1}=1}^{K}\alpha_{k_{1}k}^{(i)}\right).$$
(43)

Specifically, for i = 1, we have $\alpha_{k_1k}^{(1)} = \nu_{k_1}\nu_k$, if $k_1 \neq k$. And $\alpha_{k_1k}^{(1)} = \xi\nu_k$, if $k_1 = k$.

Sampling ν_k and ξ : As we sample $\Pi^{(i)}$, by the definition of Dirichlet-multinomial distribution, we obtain $I^{(1)} + a^{(2)} = I^{(1)} + a^{(2)}$ as Dir Mult (*u*, *transport*)

$$(l_{1k}^{(1)} + g_{\cdot 1k}^{(2)}, \cdots, l_{Kk}^{(1)} + g_{\cdot Kk}^{(2)}) \sim \text{DirMult} (\nu_1 \nu_K, \cdots, \xi \nu_k, \cdots, \nu_K \nu_k)$$

where $l_{k_1k}^{(1)} = \sum_{t=1}^{M} l_{k_1k}^{(t)}$. In particular, with a little abuse of notation here, for Dir-Dir construction, we take $g_{\cdot k_1k}^{(2)} = h_{k_1k}^{(2)}$. We first sample

$$\begin{pmatrix} h_{k_1k}^{(1)} \mid - \end{pmatrix} \sim \begin{cases} \operatorname{CRT} \left(l_{k_1k}^{(1)} + g_{\cdot k_1k}^{(2)}, \nu_{k_1}\nu_k \right) & k_1 \neq k \\ \operatorname{CRT} \left(l_{k_1k}^{(1)} + g_{\cdot k_1k}^{(2)}, \xi\nu_k \right) & k_1 = k. \end{cases}$$

$$(44)$$

526 Then we sample

$$q_k^{(1)} \sim \text{Beta}\left(l_{\cdot k}^{(1)} + g_{\cdot k}^{(2)}, \nu_k\left(\sum_{k_1 \neq k} \nu_{k1} + \xi\right)\right).$$
 (45)

527 We further introduce

$$n_{k} = h_{kk}^{(1)} + \sum_{k_{1} \neq k} h_{k_{1}k}^{(1)} + \sum_{k_{2} \neq k} h_{kk_{2}}^{(1)} + l_{k}^{(1)}, \text{ and}$$
$$\rho_{k} = \tau_{0} \zeta^{(1)} - \ln\left(1 - q_{k}^{(1)}\right) \left(\xi + \sum_{k_{1} \neq k} \nu_{k_{1}}\right) - \sum_{k_{2} \neq k} \ln\left(1 - q_{k_{2}}^{(1)}\right) \nu_{k_{2}}.$$

528 Via Poisson-gamma conjugacy, we have

$$(\xi \mid -) \sim \operatorname{Gam}\left(\frac{\gamma_0}{K} + \sum_k h_{kk}^{(1)}, \beta - \sum_k \nu_k \ln\left(1 - q_k^{(1)}\right)\right),\tag{46}$$

$$(\nu_k \mid -) \sim \operatorname{Gam}\left(\frac{\gamma_0}{K} + n_k, \beta + \rho_k\right).$$
(47)

529

Sampling $\delta^{(t)}$ and β : Via Poisson-gamma conjugacy

$$\left(\delta^{(t)} \mid -\right) \sim \operatorname{Gam}\left(\epsilon_0 + \sum_{v=1}^V y_v^{(t)}, \epsilon_0 + \sum_{k=1}^K \theta_k^{(t)}\right).$$
(48)

531 And by gamma-gamma conjugacy, we obtain

$$(\beta \mid -) \sim \operatorname{Gam}\left(\epsilon_0 + \gamma_0, \epsilon_0 + \sum_{k=1}^{K} \nu_k\right).$$
 (49)

The full procedure of our Gibbs sampling algorithms are summarized in Algorithm 1, Algorithm 2 and Algorithm 3.

Algorithm 1 Gibbs sampling algorithm for NS-PGDS (Dir-Dir Markov construction)

Input: observed count sequence $\{y^{(t)}\}_{t=1}^{T}$, iterations \mathcal{J} . Initialize the model's rank K, hyperparameters $\gamma_0, \epsilon_0, e_0, f_0$. for iter = 1 to \mathcal{J} do Sample $\{y_{vk}^{(t)}\}_{v,k}$ via Eq.(6). Sample $\{\phi_k\}_k$ via Eq.(7). Sample $\{\delta^{(t)}\}_t$ via Eq.(48). Update $\zeta^{(t)}$ as $\zeta^{(T+1)} = 0$, $\zeta^{(t)} = \ln\left(1 + \frac{\delta^{(t)}}{\tau_0} + \zeta^{(t+1)}\right)$, $t = T, \cdots, 1$. Set $l_k^{(T+1)} = 0$. for t = T to 2 do Sample $\{l_k^{(t)}\}_k$ and $\{l_{kk_2}^{(t)}\}_{k,k_2}$ via Eq.(8) and Eq.(9) respectively. end for for t = 1 to T do Sample $\{\theta_k^{(t)}\}_k$ via Eq.(12) and Eq.(13). end for for i = 1 to I do Sample $\{\pi_k^{(i)}\}_k$ via Eq.(14) and Eq.(16), Eq.(44) and Eq.(45). Sample $\{\pi_k^{(i)}\}_k$ via Eq.(14) and Eq.(19). Sample $\{\pi_k^{(i)}\}_k$ via Eq.(21). end for Sample $\xi, \{\nu_k\}_k, \beta$ via Eq.(46), Eq.(47) and Eq.(49) respectively. end for Output posterior means: $\{\theta_k^{(1:T)}\}_k, \{\phi_k\}_k, \{\pi_k^{(i)}\}_k, \delta^{(1:T)}, \xi, \{\nu_k\}_k, \beta$.

Algorithm 2 Gibbs sampling algorithm for NS-PGDS (Dir-Gam-Dir Markov construction)

Input: observed count sequence $\{\mathbf{y}^{(t)}\}_{t=1}^{T}$, iterations \mathcal{J} . Initialize the model's rank K, hyperparameters $\gamma_0, \epsilon_0, e_0, f_0$. for iter = 1 to \mathcal{J} do Sample $\{\psi_{k}^{(t)}\}_{v,k}$ via Eq.(6). Sample $\{\phi_{k}\}_{k}$ via Eq.(7). Sample $\{\delta^{(t)}\}_{t}$ via Eq.(48). Update $\zeta^{(t)}$ as $\zeta^{(T+1)} = 0$, $\zeta^{(t)} = \ln\left(1 + \frac{\delta^{(t)}}{\tau_0} + \zeta^{(t+1)}\right)$, $t = T, \cdots, 1$. Set $l_k^{(T+1)} = 0$. for t = T to 2 do Sample $\{l_k^{(t)}\}_{k}$ and $\{l_{kk_2}^{(t)}\}_{k,k_2}$ via Eq.(8) and Eq.(9) respectively. end for for t = 1 to T do Sample $\{\theta_{k+1}^{(t)}\}_{k}$ via Eq.(12) and Eq.(13). end for for i = 1 to I do Sample $\{\varphi_{k+1}^{(i)}\}_{k}$ and $\{c_{k+1}^{(i)}\}_{k}$, via Eq.(24) and Eq.(32). Sample $\{q_{k+1}^{(i)}\}_{k}$ and $\{f_{k+1}\}_{k+1,k}$ via Eq.(22), Eq.(23), Eq.(44) and Eq.(45). Sample $\{\varphi_{k+1}\}_{k,k,n}$ and $\{g_{k+1}k_2\}_{k,k+2,k}$ via Eq.(25) and Eq.(26) respectively. Sample $\{\psi_{k+1k_2}\}_{k,k+1,k_2}$ via Eq.(28). Sample $\{\varphi_{k+1}^{(i)}\}_{k}$ via Eq.(31). Sample $\{\pi_{k}^{(i)}\}_{k}$ via Eq.(31). Sample $\{\pi_{k}^{(i)}\}_{k}$ via Eq.(44) and Eq.(49) respectively. end for Sample $\{\xi, \{\nu_k\}_{k}, \beta$ via Eq.(46), Eq.(47) and Eq.(49) respectively. end for Output posterior means: $\{\theta_{k}^{(1:T)}\}_{k}, \{\phi_{k}\}_{k}, \{\pi_{k}^{(i)}\}_{k}, \delta^{(1:T)}, \xi, \{\nu_{k}\}_{k}, \beta$.

Algorithm 3 Gibbs sampling algorithm for NS-PGDS (PR-Gam-Dir Markov construction)

Input: observed count sequence $\{y^{(t)}\}_{t=1}^{T}$, iterations \mathcal{J} . **Initialize** the model's rank K, hyperparameters $\gamma_0, \epsilon_0, e_0, f_0$. for iter = 1 to \mathcal{J} do Sample $\{y_{vk}^{(t)}\}_{v,k}$ via Eq.(6). Sample $\{\phi_k\}_k$ via Eq.(7). Sample $\{\delta^{(t)}\}_t$ via Eq.(48). Update $\zeta^{(t)}$ as $\zeta^{(T+1)} = 0, \quad \zeta^{(t)} = \ln\left(1 + \frac{\delta^{(t)}}{\tau_0} + \zeta^{(t+1)}\right), \ t = T, \cdots, 1.$ Set $l_{\cdot k}^{(T+1)} = 0$. for t = T to 2 do Sample $\{l_{k}^{(t)}\}_k$ and $\{l_{kk_2}^{(t)}\}_{k,k_2}$ via Eq.(8) and Eq.(9) respectively. for t = 1 to T do Sample $\{\theta_k^{(t)}\}_k$ via Eq.(12) and Eq.(13). end for for i = 1 to I_{i} do Sample $\{\alpha_{k_1k}^{(i)}\}_{k_1,k}$ and $\{c_k^{(i)}\}_k$ via Eq.(33) and Eq.(43). Sample $\{q_k^{(i)}\}_k$ and $\{h_{k_1k}^{(i)}\}_{k_1,k}$ via Eq.(33), Eq.(34), Eq.(44) and Eq.(45). Sample $\{g_{k_1k}\}_{k_1,k}$ via Eq.(36) and Eq.(37). Sample $\{g_{k_1k_2k}\}_{k_1,k_2,k}$ via Eq.(38). Sample $\{\gamma_k^{(i)}\}_k$ via Eq.(42). Sample $\{\psi_{kk_1k_2}\}_{k,k_1,k_2}$ via Eq.(40). Sample $\{\pi_k^{(i)}\}_k$ via Eq.(14), and Eq.(41). end for Sample ξ , $\{\nu_k\}_k$, β via Eq.(46), Eq.(47) and Eq.(49) respectively. end for Output posterior means: $\{\theta_k^{(1:T)}\}_k, \{\phi_k\}_k, \{\pi_k^{(i)}\}_k, \delta^{(1:T)}, \xi, \{\nu_k\}_k, \beta$.

534 NeurIPS Paper Checklist

535 1. Claims

- Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?
- 538 Answer: [Yes]

Justification: The main contributions of this paper are the constructions of Poisson-Gamma dynamical systems with non-stationary transition dynamics and the corresponding Gibbs sampler. The constructions can be found in sec.2 and sec.3 and the derivation of Gibbs sampler can be found in sec.4 and the appendix. The experiments have demonstrated the effectiveness and features of the proposed model.

544 Guidelines:

545

546

547

548

549

550

551

552

553

555

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
 - The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
 - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

554 2. Limitations

- Question: Does the paper discuss the limitations of the work performed by the authors?
- 556 Answer: [Yes]

557	Justification: As we discussed in the conclusion part, the length of each sub-interval is a
558	constant and is treated as a hyper-parameter of the model. In the future work, we plan to
559	design a method that can find the point of change and thus the length of each sub-interval
560	can be determined automatically.
561	Guidelines:
562	• The answer NA means that the paper has no limitation while the answer No means that
563	the paper has limitations, but those are not discussed in the paper.
564	• The authors are encouraged to create a separate "Limitations" section in their paper.
565	• The paper should point out any strong assumptions and how robust the results are to
566	violations of these assumptions (e.g., independence assumptions, noiseless settings,
567	model well-specification, asymptotic approximations only holding locally). The authors
568	should reflect on now these assumptions might be violated in practice and what the
569	The three he has the first of the height of
570	• The authors should reflect on the scope of the claims made, e.g., if the approach was
571	only tested on a few datasets or with a few runs. In general, empirical results often
572	depend on implicit assumptions, which should be articulated.
573	• The authors should reflect on the factors that influence the performance of the approach.
574	For example, a facial recognition algorithm may perform poorly when image resolution
575	is low or images are taken in low lighting. Or a speech-to-text system might not be
576	used remainly to provide closed captions for online fectures because it fails to nandle
577	The such are should discuss the computational effection as of the mean and all arithme
578	• The authors should discuss the computational efficiency of the proposed algorithms
579	and now they scale with dataset size.
580	• If applicable, the authors should discuss possible limitations of their approach to
581	address problems of privacy and fairness.
582	• While the authors might fear that complete honesty about limitations might be used by
583	reviewers as grounds for rejection, a worse outcome might be that reviewers discover
584	limitations that aren't acknowledged in the paper. The authors should use their best
585	judgment and recognize that individual actions in favor of transparency play an impor-
586	will be specifically instructed to not penalize honesty concerning limitations
588	3. Theory Assumptions and Proofs
500	Question: For each theoretical result, does the nener provide the full set of assumptions and
589 590	a complete (and correct) proof?
591	Answer: [Yes]
592 593	Justification: The derivation of the Gibbs sampler is the main theoretical part of this paper which can be found in sec.4 and the appendix.
594	Guidelines:
595	• The answer NA means that the paper does not include theoretical results.
596	• All the theorems formulas and proofs in the paper should be numbered and cross-
597	referenced.
E09	• All assumptions should be clearly stated or referenced in the statement of any theorems
596	• The proofs can either encour in the main menor on the supplemental material but if
599	• The proofs can entire appear in the main paper of the supplemental material, but in the supplemental material, the supplemental material the supplemental material.
600	proof sketch to provide intuition
001	• Inversely, any informal most provided in the case of the paper should be complemented
602	• Inversely, any informat proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material
603	by formal proofs provided in appendix of suppremental material.
604	• Theorems and Lemmas that the proof relies upon should be properly referenced.
605	4. Experimental Result Reproducionity
606	Question: Does the paper fully disclose all the information needed to reproduce the main ex-
607	perimental results of the paper to the extent that it affects the main claims and/or conclusions
608	of the paper (regardless of whether the code and data are provided or not)?
609	Answer: [Yes]

610 611	Justification: We carefully described the proposed model in sec.2 and sec.3 and the experiment details can be found in sec.6.1.
612	Guidelines:
613	• The answer NA means that the paper does not include experiments.
614	• If the paper includes experiments, a No answer to this question will not be perceived
615	well by the reviewers: Making the paper reproducible is important, regardless of
616	whether the code and data are provided or not.
617	• If the contribution is a dataset and/or model, the authors should describe the steps taken
618	to make their results reproducible or verifiable.
619	• Depending on the contribution, reproducibility can be accomplished in various ways.
620	For example, if the contribution is a novel architecture, describing the architecture fully
621	might suffice, or if the contribution is a specific model and empirical evaluation, it may
622	be necessary to either make it possible for others to replicate the model with the same
623	one good way to accomplish this, but reproducibility can also be provided via detailed
625	instructions for how to replicate the results, access to a hosted model (e.g., in the case
626	of a large language model), releasing of a model checkpoint, or other means that are
627	appropriate to the research performed.
628	• While NeurIPS does not require releasing code, the conference does require all submis-
629	sions to provide some reasonable avenue for reproducibility, which may depend on the
630	nature of the contribution. For example
631	(a) If the contribution is primarily a new algorithm, the paper should make it clear how
632	to reproduce that algorithm.
633	(b) If the contribution is primarily a new model architecture, the paper should describe
634	(a) If the contribution is a new model ($a = a$ large language model), then there should
635	(c) If the contribution is a new model for reproducing the results or a way to reproduce
637	the model (e.g., with an open-source dataset or instructions for how to construct
638	the dataset).
639	(d) We recognize that reproducibility may be tricky in some cases, in which case
640	authors are welcome to describe the particular way they provide for reproducibility.
641	In the case of closed-source models, it may be that access to the model is limited in
642	some way (e.g., to registered users), but it should be possible for other researchers
643	to have some path to reproducing or verifying the results.
644	5. Open access to data and code
645	Question: Does the paper provide open access to the data and code, with sufficient instruc-
646	tions to faithfully reproduce the main experimental results, as described in supplemental
647	material?
648	Answer: [No]
649	Justification: The authors will release the data and code as soon as possible if this paper
650	could be accepted.
651	Guidelines:
652	• The answer NA means that paper does not include experiments requiring code.
653	• Please see the NeurIPS code and data submission guidelines (https://nips.cc/
654	public/guides/CodeSubmissionPolicy) for more details.
655	• While we encourage the release of code and data, we understand that this might not be
656	including code unless this is central to the contribution (e.g. for a new open-source
658	benchmark).
659	• The instructions should contain the exact command and environment needed to run to
660	reproduce the results. See the NeurIPS code and data submission guidelines (https:
661	//nips.cc/public/guides/CodeSubmissionPolicy) for more details.
662	• The authors should provide instructions on data access and preparation, including how
663	to access the raw data, preprocessed data, intermediate data, and generated data, etc.

664 665		• The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they
666		should state which ones are omitted from the script and why.
667		• At submission time, to preserve anonymity, the authors should release anonymized
668		versions (if applicable).
669 670		• Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.
671	6.	Experimental Setting/Details
672		Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
673 674		parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?
675		Answer: [Yes]
676		Justification: The experiment details can be found in sec.6.1.
677		Guidelines:
678		• The answer NA means that the paper does not include experiments.
679		• The experimental setting should be presented in the core of the paper to a level of detail
680		that is necessary to appreciate the results and make sense of them.
681		• The full details can be provided either with the code, in appendix, or as supplemental
682	_	material.
683	7.	Experiment Statistical Significance
684 685		Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?
686		Answer: [Yes]
687		Justification: Table 1 reports the predictive performance of the proposed model and the
688 689		corresponding standard deviation. The results are computed by running the Gibbs sampling several times from different initialization.
690		Guidelines:
691		• The answer NA means that the paper does not include experiments.
692		• The authors should answer "Yes" if the results are accompanied by error bars, confi-
693		dence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper
694		• The factors of veriability that the error bars are conturing should be clearly stated (for
696		example, train/test split, initialization, random drawing of some parameter, or overall
697		run with given experimental conditions).
698		• The method for calculating the error bars should be explained (closed form formula,
699		call to a library function, bootstrap, etc.)
700		• The assumptions made should be given (e.g., Normally distributed errors).
701		• It should be clear whether the error bar is the standard deviation or the standard error
702		of the filter. • It is OK to report 1 sigma error bars, but one should state it. The authors should
703		preferably report a 2-sigma error bar than state that they have a 96% CL if the hypothesis
705		of Normality of errors is not verified.
706		• For asymmetric distributions, the authors should be careful not to show in tables or
707		figures symmetric error bars that would yield results that are out of range (e.g. negative
708		error rates).
709 710		• If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
711	8.	Experiments Compute Resources
712		Ouestion: For each experiment, does the paper provide sufficient information on the com-
713		puter resources (type of compute workers, memory, time of execution) needed to reproduce
714		the experiments?
715		Answer: [Yes]

716 717		Justification: The experiments are conducted on a server with an Intel(R) Xeon(R) CPU E5-2699Cv4 @ 2.20GHz and 64G RAM.
718		Guidelines:
710		• The answer NA means that the paper does not include experiments
700		• The paper should indicate the type of compute workers CPU or CPU internal cluster
720		or cloud provider, including relevant memory and storage.
722		• The paper should provide the amount of compute required for each of the individual
723		experimental runs as well as estimate the total compute.
724		• The paper should disclose whether the full research project required more compute
725 726		than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper)
727	9.	Code Of Ethics
728		Question: Does the research conducted in the paper conform, in every respect, with the
729		NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?
730		Answer: [Yes]
731 732		Justification: We have checked the NeurIPS Code of Ethics and make sure this work is with the NeurIPS Code of Ethics.
733		Guidelines:
		• The answer NA means that the outhers have not reviewed the NeurIDS Code of Ethics
/34		• The answer IVA means that the authors have not reviewed the recurrence of European authors are such as a straight of the second authors are straighted authors are straighted as a second author are straighted as a second at the
735		• If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics
730		• The authors should make sure to preserve anonymity (e.g. if there is a special consid-
738		eration due to laws or regulations in their jurisdiction)
700	10	Broader Impacts
739	10.	Droader impacts
740 741		societal impacts of the work performed?
742		Answer: [Yes]
743		Justification: For positive societal impacts, we have discussed in conclusion section for the
744		potential application for textual analysis and social networks. And the authors think this
745		work does not have potential negative societal impacts.
746		Guidelines:
747		 The answer NA means that there is no societal impact of the work performed.
748		• If the authors answer NA or No, they should explain why their work has no societal
749		impact or why the paper does not address societal impact.
750		• Examples of negative societal impacts include potential malicious or unintended uses
751		(e.g., disinformation, generating take profiles, surveillance), fairness considerations
752		(e.g., deployment of technologies that could make decisions that unfairly impact specific
/53		groups), privacy considerations, and security considerations.
754		• The conference expects that many papers will be foundational research and not fied to particular applications, let along deployments. However, if there is a direct path to
/55		any negative applications, the authors should point it out. For example, it is legitimate
750		to point out that an improvement in the quality of generative models could be used to
758		generate deepfakes for disinformation. On the other hand, it is not needed to point out
759		that a generic algorithm for optimizing neural networks could enable people to train
760		models that generate Deepfakes faster.
761		• The authors should consider possible harms that could arise when the technology is
762		being used as intended and functioning correctly, harms that could arise when the
763		technology is being used as intended but gives incorrect results, and harms following
764		from (intentional or unintentional) misuse of the technology.
765		• If there are negative societal impacts, the authors could also discuss possible mitigation
766		strategies (e.g., gated release of models, providing defenses in addition to attacks,
767		mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
768		reedback over time, improving the efficiency and accessibility of ML).

769	11.	Safeguards
770		Question: Does the paper describe safeguards that have been put in place for responsible
771		release of data or models that have a high risk for misuse (e.g., pretrained language models,
772		image generators, or scraped datasets)?
773		Answer: [NA]
774		Justification: This work poses no such risks.
775		Guidelines:
776		• The answer NA means that the paper poses no such risks.
777		• Released models that have a high risk for misuse or dual-use should be released with
778		necessary safeguards to allow for controlled use of the model, for example by requiring
779		that users adhere to usage guidelines or restrictions to access the model or implementing
780		safety filters.
781 782		• Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
783		• We recognize that providing effective safeguards is challenging, and many papers do
784		not require this, but we encourage authors to take this into account and make a best
785		faith effort.
786	12.	Licenses for existing assets
787		Question: Are the creators or original owners of assets (e.g., code, data, models), used in
788		the paper, properly credited and are the license and terms of use explicitly mentioned and
789		properly respected?
790		Answer: [NA]
791		Justification: This paper does not use existing assets.
792		Guidelines:
793		• The answer NA means that the paper does not use existing assets.
794		• The authors should cite the original paper that produced the code package or dataset.
795 796		• The authors should state which version of the asset is used and, if possible, include a URL
797		• The name of the license (e.g. CC-BY 4.0) should be included for each asset
798		• For scraped data from a particular source (e.g. website) the convright and terms of
799		service of that source should be provided.
800		• If assets are released, the license, copyright information, and terms of use in the
801		package should be provided. For popular datasets, paperswithcode.com/datasets
802		has curated licenses for some datasets. Their licensing guide can help determine the
803		Experience of a dataset.
804 805		• For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
806		• If this information is not available online, the authors are encouraged to reach out to
807		the asset's creators.
808	13.	New Assets
809		Question: Are new assets introduced in the paper well documented and is the documentation
810		provided alongside the assets?
811		Answer: [NA]
812		Justification: This paper does not release new assets.
813		Guidelines:
814		• The answer NA means that the paper does not release new assets.
815		• Researchers should communicate the details of the dataset/code/model as part of their
816		submissions via structured templates. This includes details about training, license,
817		limitations, etc.
818		• The paper should discuss whether and how consent was obtained from people whose
819		asset is used.

820 821	• At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
822 14	4. Crowdsourcing and Research with Human Subjects
823 824 825	Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?
826	Answer: [NA]
827	Justification: This paper does not involve crowdsourcing nor research with human subjects.
828	Guidelines:
829 830	• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
831 832 833	• Including this information in the supplemental material is fine, but if the main contribu- tion of the paper involves human subjects, then as much detail as possible should be included in the main paper.
834	• According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
835 836	or other labor should be paid at least the minimum wage in the country of the data collector.
837 1.	5. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
838	Subjects
839 840 841 842	Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?
843	Answer: [NA]
844	Justification: This paper does not involve crowdsourcing nor research with human subjects.
845	Guidelines:
846 847	• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
848 849 850	• Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
851 852 853	• We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
854 855	• For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.