
Poisson-Randomized Gamma Dynamical Systems

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 This paper presents the Poisson-randomized gamma dynamical system (PrGDS),
2 a model for sequentially-observed count tensors that encodes a strong inductive
3 bias towards sparsity and burstiness. The PrGDS is based on a new motif in
4 Bayesian latent variable modeling: an alternating series of discrete Poisson and
5 continuous gamma latent states. This motif is widely applicable and analytically
6 tractable, yielding closed-form complete conditionals for all variables by way of
7 the Bessel distribution and a novel distribution that we call the shifted confluent
8 hypergeometric distribution. We draw connections to closely-related models
9 and compare the PrGDS to them in studies of real-world count data of text,
10 international events, and neural spike trains. We find that a sparse variant of
11 the PrGDS—which allows continuous latent states to take values of exactly
12 zero—often obtains the lowest smoothing and forecasting perplexity of all models
13 and is uniquely capable of inferring latent structure that is highly localized in time.

14 1 Introduction

15 Political scientists regularly analyze counts of the number of times country i took action a towards
16 country j during time step t [1]. Such data exhibits “complex dependence structures” [2] like
17 coalitions of countries and bursty temporal dynamics. These dependence structures violate the
18 independence assumptions of traditional regression methods that political scientists have traditionally
19 used to test theories of international relations [3, 4, 5]. Political scientists have thus advocated for
20 using latent variable models to infer unobserved structure as a way of controlling for it [6]. The latter
21 approach motivates interpretable yet expressive models, capable of capturing a variety of complex
22 latent structures. Event data sets can be represented as a sequence of count tensors $\mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(T)}$
23 each of which contains the $V \times V \times A$ event counts for that time step for every combination of V sender
24 countries, V receivers, and A action types. Recent work applies tensor decomposition methods to
25 event data sets [7, 8, 9, 10, 11], which infers interpretable coalition structure among countries and topic
26 structure among actions. Like most tensor decomposition methods though, these methods assume
27 that the sequence of tensors is exchangeable and cannot capture the temporal structure in the data.

28 Sequentially observed count tensors present unique statistical challenges because they tend to be
29 *bursty* [12], *high-dimensional*, and *sparse* [13, 14]. There are few models that are tailored to both
30 the challenging properties count time-series and count tensors. In recent years, Poisson factorization
31 has emerged as a framework for modeling sparse count matrices [15, 16, 17, 18, 19, 20] and tensors
32 [13, 21, 9]. While tensor decomposition methods generally scale with the size the tensor, many
33 Poisson factorization models yield inference algorithms that scale linearly with only the non-zeros.
34 This property allows researchers to efficiently explore latent structure in massive tensors, provided
35 they are sparse. However, this property is unique to a subset of Poisson factorization models that only
36 use non-negative prior distributions, which are difficult to chain in state-space models for time series.
37 Hierarchical compositions of non-negative priors—notably, the gamma and Dirichlet—typically
38 introduce non-conjugate dependencies that require innovative posterior schemes.

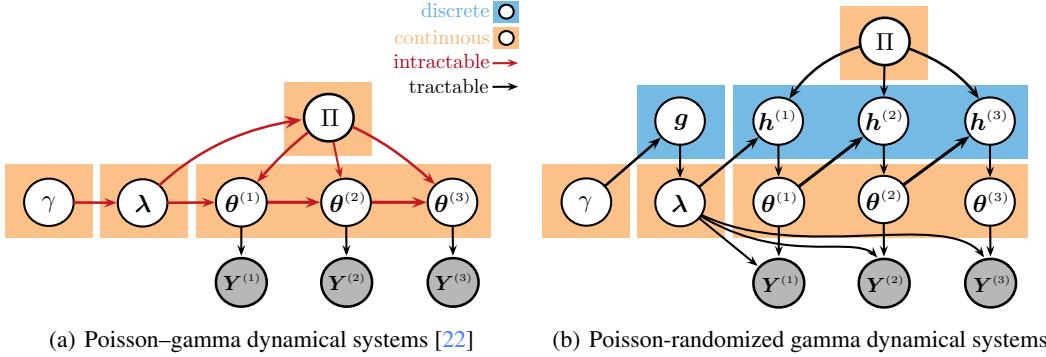


Figure 1: *Left:* The PGDS imposes dependencies directly between the continuous variables that do not yield closed-form conditional distributions. *Right:* The PrGDS (this paper) breaks the intractable dependencies with discrete Poisson variables—doing so yields closed-form conditionals for all variables without any data augmentation.

This paper seeks to fill a gap in the literature between Poisson factorization models that are *tractable*—i.e., yielding closed-form complete conditionals that make approximate inference easy to derive—and those that are *expressive*—i.e., capable of capturing a variety of complex dependence structures. To do so, we introduce alternating chains of discrete Poisson and continuous gamma latent states, a new modeling motif that is analytically convenient and computationally tractable. We rely on this motif to construct the Poisson-randomized gamma dynamical system (PrGDS), a model for sequentially observed count tensors that is tractable, expressive, and efficient. The PrGDS is closely related to the Poisson–gamma dynamical system (PGDS) [22], a recently introduced model for dynamic count matrices, that is based on non-conjugate chains of gamma-distributed states. These chains are intractable—thus, posterior inference in the PGDS relies on sophisticated data augmentation schemes that are cumbersome to derive and impose unnatural restrictions on the priors over other variables. The PrGDS instead introduces intermediate Poisson states that break the intractable dependency between the gamma states (see Fig. 1). While this construction is only *semi*-conjugate, it is tractable, yielding closed-form complete conditionals for the Poisson states by way of the little-known Bessel distribution [23] and a novel discrete distribution that we derive and call the *shifted confluent hypergeometric (SCH) distribution*.

We study the inductive bias of the PrGDS by comparing its smoothing and forecasting ability to the PGDS and two other baselines on a range of real-world count matrices and tensors of text, international events, and neural spike data. We find that the PrGDS often obtains lower smoothing and forecasting perplexity than the PrGDS and related baselines. The PrGDS under a specific hyperparameter settings permits the continuous states to take values of *exactly* zero thus encoding a unique inductive bias tailored to sparsity and burstiness. We find that this variant, in particular, often obtains the lowest perplexity of all models. We also find that the sparse PrGDS is representing of inferring a qualitatively broader range of latent structure—specifically, bursty latent structure that is highly localized in time.

2 Poisson-randomized gamma dynamical systems (PrGDS)

Notation. Consider a data set of sequentially observed tensors $\mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(T)}$. An entry $y_{\mathbf{i}}^{(t)} \in \{0, 1, 2, \dots\}$ in the t^{th} tensor is subscripted by a multi-index $\mathbf{i} \equiv (i_1, \dots, i_M)$ which indexes into the M modes of the tensor. As an example, international event counts $y_{i \xrightarrow{a_j} j}^{(t)}$ collectively form a sequence of 3-mode count tensors where each multi-index corresponds to a unique combination of sender, receiver, and action type—e.g., $\mathbf{i} = (i, j, a)$.

Generative process. The PrGDS is a form of canonical polyadic decomposition [24] that models $y_{\mathbf{i}}$ as

$$y_{\mathbf{i}}^{(t)} \sim \text{Pois}\left(\rho^{(t)} \sum_{k=1}^K \lambda_k \theta_k^{(t)} \prod_{m=1}^M \phi_{k i_m}^{(m)}\right). \quad (1)$$

Here $\theta_k^{(t)}$ represents the activation of the k^{th} component at time step t . Each component describes a dependence structure in the data by way of a factor vector $\phi_k^{(m)}$ for each mode m . For international events data, the first factor vector $\phi_k^{(1)} = (\phi_{k1}^{(1)}, \dots, \phi_{kV}^{(1)})$ would describe the rate at which each of the V countries acts as a sender in the k^{th} component while the second $\phi_k^{(2)}$ would describe the rate

73 at which each acts as a receiver. The weights λ_k and $\rho^{(t)}$ represent the overall scale of component k
 74 and time step t . The PrGDS is called *stationary* if $\rho^{(t)} = \rho$. We posit the following conjugate priors,

$$\rho^{(t)} \sim \text{Gam}(a_0, b_0) \quad \text{and} \quad \phi_k^{(m)} \sim \text{Dir}(a_0, \dots, a_0). \quad (2)$$

75 The PrGDS is characterized by an alternating series of discrete and continuous latent states. The
 76 *continuous latent states* $\theta_k^{(t)}$ evolve via intermediate discrete states $h_k^{(t)}$ from $t = 1, \dots, T$ as

$$\theta_k^{(t)} \sim \text{Gam}(\epsilon_0^{(\theta)} + h_k^{(t)}, \tau) \quad \text{and} \quad h_k^{(t)} \sim \text{Pois}\left(\tau \sum_{k_2=1}^K \pi_{kk_2} \theta_{k_2}^{(t-1)}\right), \quad (3)$$

77 where for $t = 0$ we define $\theta_k^{(0)} = \lambda_k$ to be the per-component weight that also appears in Eq. (1).
 78 The PrGDS assumes $\theta_k^{(t)}$ is conditionally gamma distributed with *rate* τ and shape equal to a latent
 79 count $h_k^{(t)}$ plus hyperparameter $\epsilon_0^{(\theta)} \geq 0$. We adopt the convention that a gamma random variable
 80 is zero, almost surely, if its shape parameter is zero—thus, setting $\epsilon_0^{(\theta)} = 0$ defines the *sparse PrGDS*
 81 wherein the continuous states are exactly zero $\theta_k^{(t)} = 0$ if $h_k^{(t)} = 0$. The *transition weights* π_{kk_2} in the
 82 Poisson rate of $h_k^{(t)}$ represent how strongly each component k_2 excites component k at the subsequent
 83 time step. We view these weights collectively as a $K \times K$ transition matrix Π and impose Dirichlet
 84 priors over its columns. We also place a gamma prior over the *concentration parameter* τ which
 85 is conjugate to both the gamma and Poisson distributions it appears in:

$$\tau \sim \text{Gam}(\alpha_0, \alpha_0) \quad \text{and} \quad \pi_k \sim \text{Dir}(a_0, \dots, a_0) \quad \text{such that } \sum_{k_1}^K \pi_{k_1 k} = 1. \quad (4)$$

86 For the per-component weights λ_k , we place a hierarchical prior with a similar flavor to Eq. (3):

$$\lambda_k \sim \text{Gam}\left(\frac{\epsilon_0^{(\lambda)}}{K} + g_k, \beta\right) \quad \text{and} \quad g_k \sim \text{Pois}\left(\frac{\gamma}{K}\right), \quad (5)$$

87 where $\epsilon_0^{(\lambda)}$ is a hyperparameter analogous to $\epsilon_0^{(\theta)}$. The following gamma priors are then both conjugate:

$$\gamma \sim \text{Gam}(a_0, b_0) \quad \text{and} \quad \beta \sim \text{Gam}(\alpha_0, \alpha_0). \quad (6)$$

88 **Properties.** Both $\epsilon_0^{(\lambda)}$ and γ are divided by the number of components K in Eq. (5)—as the number
 89 of components grows $K \rightarrow \infty$, the expected sum of the weights thus remains finite and fixed:

$$\sum_{k=1}^{\infty} \mathbb{E}[\lambda_k] = \sum_{k=1}^{\infty} \left(\frac{\epsilon_0^{(\lambda)}}{K} + \mathbb{E}[g_k]\right) \beta^{-1} = \sum_{k=1}^{\infty} \left(\frac{\epsilon_0^{(\lambda)}}{K} + \frac{\gamma}{K}\right) \beta^{-1} = (\epsilon_0^{(\lambda)} + \gamma) \beta^{-1}. \quad (7)$$

90 Thus, this prior encodes an inductive bias towards small values of λ_k and may be interpreted as the
 91 finite truncation of a novel Bayesian nonparametric process. A small value of λ_k shrinks the Poisson
 92 rates of both the data $y_i^{(t)}$ and the first discrete latent state $h_k^{(0)}$ —this prior encourages the model to
 93 only infer components that are both predictive of the data and useful for fitting the latent dynamics.

94 The marginal expectation of the state vector $\boldsymbol{\theta}^{(t)}$ takes the canonical form of linear dynamical systems,

$$\mathbb{E}[\boldsymbol{\theta}^{(t)} | \boldsymbol{\theta}^{(t-1)}] = \mathbb{E}[\mathbb{E}[\boldsymbol{\theta}^{(t)} | \mathbf{h}^{(t-1)}]] = \epsilon_0^{(\theta)} \tau^{-1} + \Pi \boldsymbol{\theta}^{(t-1)}, \quad (8)$$

95 since by iterated expectation $\mathbb{E}[\theta_k^{(t)}] = (\epsilon_0^{(\theta)} + \mathbb{E}[h_k^{(t)}]) \tau^{-1} = (\epsilon_0^{(\theta)} + \tau \sum_{k_2=1}^K \pi_{kk_2} \theta_{k_2}^{(t-1)}) \tau^{-1}$. The
 96 *concentration parameter* τ appears in both the Poisson and gamma distributions in Eq. (3). It
 97 contributes to the variance of the process while canceling out of the expectation, except for the
 98 additive term $\epsilon_0^{(\theta)} \tau^{-1}$ which vanishes when $\epsilon_0^{(\theta)} = 0$.

99 More generally, we can marginalize out all of the discrete latent states $h_k^{(t)}$ to obtain a purely continuo
 100 us dynamical system in terms of the *randomized gamma distribution of the first type* (RG1) [23, 25],

$$\theta_k^{(t)} \sim \text{RG1}\left(\epsilon_0^{(\theta)}, \tau \sum_{k_2=1}^K \pi_{kk_2} \theta_{k_2}^{(t-1)}, \tau\right), \quad (9)$$

101 when $\epsilon_0^{(\theta)} > 0$ and in terms of a limiting form of the RG1 when $\epsilon_0^{(\theta)} = 0$. We describe the RG1 in Fig. 2.

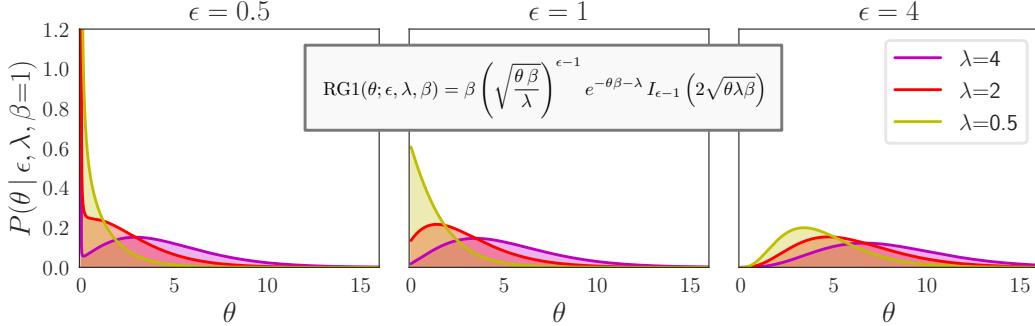


Figure 2: The randomized gamma distribution of the first type (RG1) [23, 25] has support $\theta > 0$ and is defined by three parameters $\epsilon, \lambda, \beta > 0$. Its PDF is given in the plot where $I_\nu(a)$ is the modified Bessel function of the first kind [26]. When $\epsilon < 1$ (Left) the RG1 resembles a soft “spike-and-slab” while when $\epsilon \geq 1$ (Middle and Right) it resembles a more-dispersed form of the gamma distribution. A limiting case of the RG1 when $\epsilon \rightarrow 0$ is the Poisson-randomized gamma distribution [27] which includes zeros in its support $\theta \geq 0$.

102 3 Related work

103 The PrGDS closely relates to the Poisson–gamma dynamical system (PGDS) [22] wherein

$$\theta_k^{(t)} \sim \text{Gam}\left(\tau \sum_{k_2=1}^K \pi_{kk_2} \theta_{k_2}^{(t-1)}, \tau\right) \text{ such that } \mathbb{E}[\theta^{(t)} | \theta^{(t-1)}] = \Pi \theta^{(t-1)}. \quad (10)$$

104 The PGDS posits non-conjugate dependence between gamma-distributed states. The complete
105 conditional $P(\theta_k^{(t)} | -)$ is not available in closed form under the PDGS and posterior inference relies on
106 a sophisticated data augmentation scheme. The PrGDS instead introduces intermediate Poisson states
107 $h_k^{(t)}$ that break the intractable dependencies between the gamma states. The Poisson is not a conjugate
108 prior to the gamma—however, this construction is still tractable, and yields a closed-form conditional
109 $P(h_k^{(t)} | -)$, as we’ll show in § 4. The PGDS is limited by the data augmentation scheme it requires
110 for posterior inference—specifically, this augmentation scheme does not allow any per-component
111 weights λ_k to appear in the Poisson rate of $y_i^{(t)}$ in Eq. (1). To encourage parsimony, the PGDS instead
112 draws weights $\lambda_k \sim \text{Gam}(\frac{\gamma}{K}, \beta)$ and uses them to shrink the transition matrix Π . This introduces
113 more intractable dependencies that necessitate a different augmentation scheme for inference. These
114 augmentation schemes additionally impose restrictions that the factors $\phi_k^{(m)}$ and columns π_k of the
115 transition matrix are Dirichlet distributed—while we adopt those same assumptions in this paper, the
116 PrGDS is not bound to them. We provide a graphical comparison of PGDS to PrGDS in Fig. 1.

117 The PGDS and its “deep” variants [28, 29], generalize gamma process dynamic Poisson factor
118 analysis (GP-DPFA) [30] which assumes a simple random walk $\theta_k^{(t)} \sim \text{Gam}(\theta_k^{(t-1)}, c^{(t)})$; see also a
119 related model [31]. These models belong to a line of work exploring the application of the “augment-
120 and-conquer” data augmentation scheme [32] to perform inference in hierarchies of gamma variables
121 chained via their shape and linked to Poisson observations—beyond time-series models, this construc-
122 tion can be used to build belief networks [33]. An alternative approach is to chain gamma variables
123 via their rate—e.g., $\theta^{(t)} \sim \text{Gam}(a, \theta^{(t-1)})$. This motif is conjugate and tractable and has been applied
124 to time-series models [34, 35, 36] as well as deep belief networks [37]. However, the rate contributes
125 quadratically to the variance of the gamma distribution and these chains may be highly volatile.

126 More broadly, gamma shape and rate chains are examples of non-negative chains. Such chains
127 are particularly motivated in the context of Poisson factorization, which is particularly efficient
128 when only non-negative prior distributions are used. In general, Poisson factorization assumes
129 that each observed count is drawn $y_i \sim \text{Pois}(\mu_i)$ with latent rate μ_i defined to be some function of
130 model parameters. When the rate is linear—i.e., $\mu_i = \sum_k \mu_{ik}$ —Poisson factorization yields a *latent*
131 *source representation* [16, 18] wherein we define $y_i \triangleq \sum_k y_{ik}$ to be the sum of latent sources, each
132 of which is drawn $y_{ik} \sim \text{Pois}(\mu_{ik})$. Conditioning on the latent sources often induces conditional
133 independences between the latent variables and parameters that facilitates closed-form, efficient,
134 and parallelizable posterior inference—thus, the first step in either MCMC or variational inference
135 is to update the latent sources from their complete conditionals, which is multinomial [38],

$$((y_{i1}, \dots, y_{iK}) | -) \sim \text{Multinom}(y_i, (\mu_{i1}, \dots, \mu_{iK})), \quad (11)$$

136 where we leave implicit the normalization of the non-negative rates μ_{ik} into a probability vector.
 137 When the observed count is zero $y_i = 0$, the sources are zero $y_{ik} = 0$, almost surely, and no compu-
 138 tation is required to update them. Thus, any Poisson factorization model that admits a latent source
 139 representation scales linearly with only the non-zero counts in the data. This property is indispensable
 140 when modeling count tensors which typically contain exponentially more entries than non-zeros [39].
 141 Since the latent source representation is only available when the rate μ_i is a (multi)linear function of
 142 parameters and since the rate must be non-negative, by definition of the Poisson distribution, efficient
 143 Poisson factorization is only compatible with non-negative priors over parameters. Modeling time
 144 series and other complex dependence structures with efficient Poisson factorization often requires
 145 developing novel motifs that notably exclude Gaussian priors which researchers have traditionally
 146 relied on for their analytic convenience and tractability. The Poisson LDS, for instance, [40] links the
 147 widely-used Gaussian linear dynamical system [41, 42] to Poisson observations via an exponential
 148 link function $\mu_i = \exp(\sum_k \dots)$. This is one instance of the generalized linear model (GLM) [43]
 149 approach that relies on a non-linear link function and thus does not yield a latent source representation.
 150 Another approach is to use log-normal priors, as in Dynamic Poisson Factorization [44]—while this
 151 approach satisfies the non-negative constraint, the log-normal is not conjugate to the Poisson and
 152 does not closed-form conditionals. We also note a long tradition of autoregressive models for count
 153 time series including VAR models [45] and those based on Hawkes processes [46, 47, 48]. This
 154 approach avoids the challenge of constructing tractable state-space models from non-negative priors
 155 by modeling temporal correlation directly between the count observations. For high-dimensional
 156 data, such as sequentially-observed tensors, an autoregressive approach is often untenable.

157 4 Posterior inference

158 The complete conditionals for all latent variables in the PrGDS are immediately available in closed
 159 form without any data augmentation. Iteratively re-sampling each variable from its conditional
 160 constitutes a Gibbs sampling algorithm. We provide conditionals for the latent variables with
 161 non-standard priors here and relegate the rest to the Appendix. The PrGDS is based on a new
 162 modeling motif—we first introduce it in its general form, derive its complete conditionals, and then
 163 apply these identities to the PrGDS.

164 4.1 Poisson–gamma–Poisson recursions

165 Consider the following model of m involving latent variables θ and h and fixed $c_1, c_2, c_3, \epsilon_0^{(\theta)} > 0$:

$$m \sim \text{Pois}(\theta c_3), \quad \theta \sim \text{Gam}(\epsilon_0^{(\theta)} + h, c_2), \quad \text{and } h \sim \text{Pois}(c_1). \quad (12)$$

166 This model is *semi*-conjugate. The gamma prior of θ is conjugate to the Poisson and its posterior is

$$(\theta | -) \sim \text{Gam}(\epsilon_0^{(\theta)} + h + m, c_2 + c_3). \quad (13)$$

167 The Poisson prior over h is not conjugate to the gamma—however the conditional posterior of h is still available in closed form by way of the Bessel distribution [23] which we define in Fig. 3(a):

$$(h | -) \sim \text{Bes}(\epsilon_0^{(\theta)} - 1, 2\sqrt{\theta c_2 c_1}). \quad (14)$$

168 The Bessel distribution can be sampled efficiently [49]; we will release our Cython implementation.
 169 Provided that $\epsilon_0^{(\theta)} > 0$, sampling θ and h iteratively from Eqs. (13) and (14) constitutes a valid
 170 Markov chain for posterior inference. When $\epsilon_0^{(\theta)} = 0$ though, $\theta = 0$, almost surely, if $h = 0$, and vice
 171 versa—thus, this Markov chain has an absorbing condition at $h = 0$ and violates detailed balance. In
 172 this case, we must therefore sample h with θ marginalized out—towards that end, we give Theorem 1.

173 **Theorem 1:** *The incomplete conditional $P(h | \epsilon_0^{(\theta)} = 0, - \setminus \theta) \triangleq \int P(h, \theta | \epsilon_0^{(\theta)} = 0, -) d\theta$ is*

$$(h | - \setminus \theta) \sim \begin{cases} \text{Pois}\left(\frac{c_1 c_2}{c_3 + c_2}\right) & \text{if } m = 0 \\ \text{SCH}\left(m, \frac{c_1 c_2}{c_3 + c_2}\right) & \text{otherwise} \end{cases} \quad (15)$$

174 where SCH is a discrete distribution we call the shifted confluent hypergeometric distribution. We de-
 175 scribe the SCH in Fig. 3(b) and provide further details about it in the Appendix including the derivation
 176 of its PMF, PGF, and mode, along with details of how we sample from it and the proof for Theorem 1.

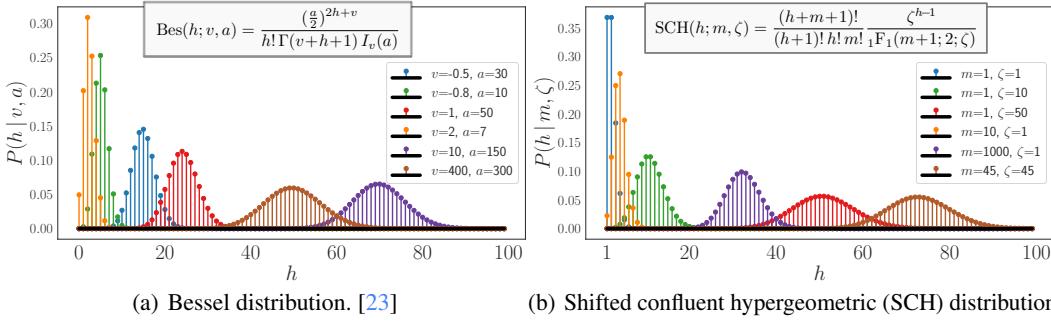


Figure 3: Two discrete distributions that arise as posteriors in Poisson–gamma–Poisson recursions.

177 4.2 Closed-form complete conditionals for the PrGDS

178 The PrGDS yields a latent source representation (see Eq. (11))—posterior inference begins with

$$((y_{ik}^{(t)})_{k=1}^K | -) \sim \text{Multinom}\left(y_i^{(t)}, (\lambda_k \theta_k^{(t)} \prod_{m=1}^M \phi_{ki_m}^{(m)})_{k=1}^K\right). \quad (16)$$

We may similarly represent $h_k^{(t)}$ under its latent source representation as $h_k^{(t)} \equiv h_{k \cdot}^{(t)} = \sum_{k_2=1}^K h_{kk_2}^{(t)}$ where $h_{kk_2}^{(t)} \sim \text{Pois}(\tau \pi_{kk_2} \theta_{k_2}^{(t-1)})$. When useful, we use dot-notation (“.”) to denote summing over an axis—in this case $h_{k \cdot}^{(t)}$ denotes the sum of the k^{th} row of the $K \times K$ matrix of latent counts $h_{kk_2}^{(t)}$.

179 The complete conditional of the k^{th} row of counts, when conditioned on their sum $h_{k \cdot}^{(t)}$, is

$$((h_{k \cdot}^{(t)})_{k_2=1}^K | -) \sim \text{Multinom}\left(h_{k \cdot}^{(t)}, (\pi_{kk_2} \theta_{k_2}^{(t-1)})_{k_2=1}^K\right). \quad (17)$$

180 To derive the conditional for $\theta_k^{(t)}$ we aggregate all Poisson variables that depend on it. By Poisson additivity, the column sum $h_{\cdot k}^{(t+1)} = \sum_{k_1=1}^K h_{k_1 k}^{(t+1)}$ is distributed $h_{\cdot k}^{(t+1)} \sim \text{Pois}(\theta_k^{(t)} \tau \pi_{\cdot k})$ and similarly $y_{\cdot k}^{(t)}$ is distributed $y_{\cdot k}^{(t)} \sim \text{Pois}(\theta_k^{(t)} \rho^{(t)} \lambda_k \prod_{m=1}^M \phi_{k \cdot}^{(m)})$. The count $m_k^{(t)} \triangleq h_{\cdot k}^{(t+1)} + y_{\cdot k}^{(t)}$ then isolates all dependence on $\theta_k^{(t)}$ and is also Poisson distributed. By gamma–Poisson conjugacy, the conditional of $\theta_k^{(t)}$ is

$$(\theta_k^{(t)} | -) \sim \text{Gam}(\epsilon_0^{(\theta)} + h_{k \cdot}^{(t)} + m_k^{(t)}, \tau + \tau \pi_{\cdot k} + \rho^{(t)} \lambda_k \prod_{m=1}^M \phi_{k \cdot}^{(m)}). \quad (18)$$

184 When $\epsilon_0^{(\theta)} > 0$, we may apply the identity in Eq. (14) and sample $h_{k \cdot}^{(t)}$ from its complete conditional:

$$(h_{k \cdot}^{(t)} | -) \sim \text{Bessel}\left(\epsilon_0^{(\theta)} - 1, 2\sqrt{\theta_k^{(t)} \tau^2 \sum_{k_2=1}^K \pi_{kk_2} \theta_{k_2}^{(t-1)}}\right). \quad (19)$$

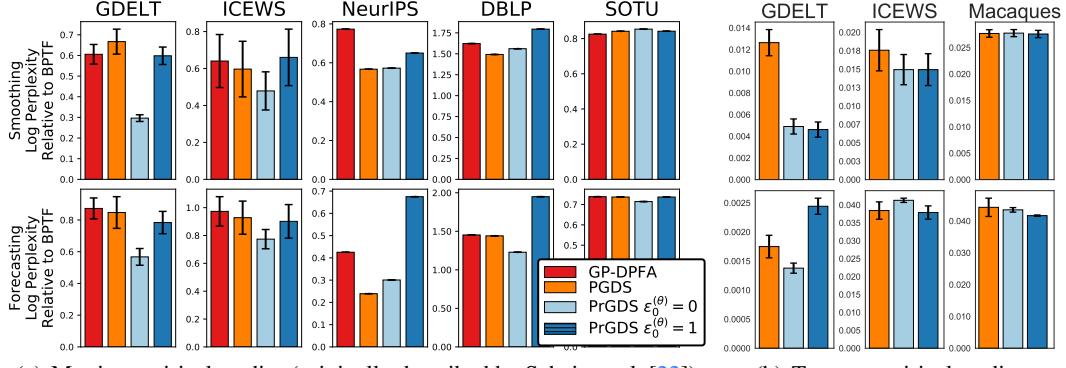
185 When $\epsilon_0^{(\theta)} = 0$, we instead apply Theorem 1 to sample $h_{k \cdot}^{(t)}$ where $m_k^{(t)}$ is analogous to m in Eq. (15):

$$(h_{k \cdot}^{(t)} | - \setminus \theta_k^{(t)}) \sim \begin{cases} \text{Pois}(\zeta_k^{(t)}) & \text{if } m_k^{(t)} = 0 \\ \text{SCH}(m_k^{(t)}, \zeta_k^{(t)}) & \text{otherwise} \end{cases} \quad \text{where } \zeta_k^{(t)} \triangleq \frac{\tau^2 \sum_{k_2=1}^K \pi_{kk_2} \theta_{k_2}^{(t-1)}}{\tau + \tau \pi_{\cdot k} + \rho^{(t)} \lambda_k \prod_{m=1}^M \phi_{k \cdot}^{(m)}}. \quad (20)$$

186 The conditionals for λ_k and g_k follow from applying the same Poisson–gamma–Poisson identities while those for γ , β , $\phi_k^{(m)}$, π_k , and τ all follow from conjugacy. We provide them all in the Appendix.

188 5 Empirical study

189 We’ve seen how the Poisson–gamma–Poisson motif of the PrGDS (see § 4.1) yields a more tractable
190 (see Fig. 1) and more flexible (see § 3) model family than previous work. This motif also encodes
191 a unique inductive bias that we empirically test by comparing the PrGDS to the Poisson–gamma
192 dynamical system (PGDS) [22]. The PGDS is the pure gamma analog to the PrGDS, as we
193 see by comparing Eqs. (9) and (10)—comparing these models thus isolates the impact of the
194 Poisson–gamma–Poisson motif. The PGDS was only previously introduced to model a $T \times V$ matrix
195 Y of sequentially observed V -dimensional vectors $y^{(1)}, \dots, y^{(T)}$. To compare to it, we generalize the
196 PGDS to M -mode tensors. We have provided our Cython implementation of the generalized PGDS
197 (in addition to the PrGDS) and derive its complete conditionals in the Appendix. We also compare
198 the PrGDS variant with $\epsilon_0^{(\theta)} = 1$ to the one with $\epsilon_0^{(\theta)} = 0$, which permits sparse activations $\theta_k^{(t)} = 0$.



(a) Matrix empirical studies (originally described by Schein et al. [22]). (b) Tensor empirical studies.

Figure 4: The smoothing (top row) and forecasting (bottom row) performance of each model is measured by log-perplexity—where *lower is better*—divided by the log-perplexity of a non-dynamic baseline, BPTF [9].

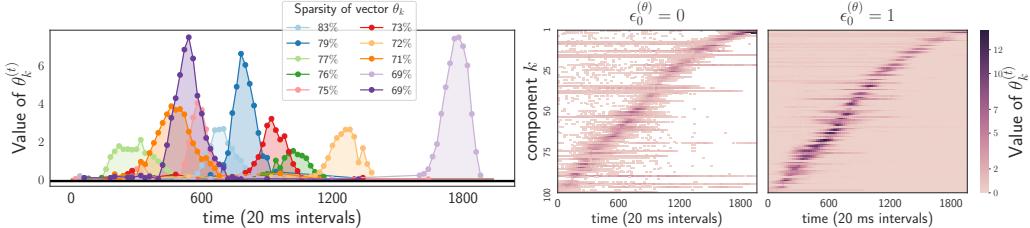
199 **Setup.** All empirical studies have the following setup: for some data set $\mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(T)}$, all of
200 the counts $\mathbf{Y}^{(t)}$ at randomly selected time steps t are masked. The last two time steps are also
201 always masked. Each model is fit to the masked data using independent chains of MCMC that
202 impute the missing data each at iteration and ultimately return S posterior samples of the latent
203 variables and parameters. The samples are then used to compute the expectations $\mu_i^{(t)}$ (as defined
204 in Eq. (1)) of the heldout counts. We distinguish the task of predicting counts in the final time steps
205 when subsequent observed data is unavailable (*forecasting*) from the task of predicting intermediate
206 time steps (*smoothing*). To assess predictive performance we compute log-perplexity as defined—
207 $\log \text{Perp}(\Delta) = -\frac{1}{|\Delta|} \sum_{(t,i) \in \Delta} \log \left[\frac{1}{S} \sum_{s=1}^S \text{Pois}(y_i^{(t)}; \mu_{i,s}^{(t)}) \right]$ —where Δ is the set of multi-indices
208 of the heldout counts and $\mu_{i,s}^{(t)}$ is the expectation of the heldout count computed from the s^{th} sample.
209 In all studies, we fit a simple non-dynamic baseline—i.e., Bayesian Poisson tensor factorization
210 (BPTF) [9]—that assumes the data at different time slices are i.i.d. $y_i^{(t)} \sim \text{Pois}(\mu_i)$. This model thus
211 fits μ_i from training data and predicts it for all heldout time slices. For the dynamic models, we then
212 report their improvement over non-dynamic BPTF by dividing their log-perplexity by BPTF’s.

213 **Matrices.** We first replicated the empirical studies on $T \times V$ dynamic *matrices* (i.e., 2-mode tensors)
214 reported by Schein et al. (2016) [22]. These studies followed the setup described above and compared
215 the PGDS to GP-DPFA [30], a simple dynamic baseline that we describe in § 3. The matrices in
216 these studies are based on three text corpora—NeurIPS papers [50], DBLP abstracts [51], and State
217 of the Union (SOTU) speeches [52]—where $y_v^{(t)}$ is the number of times word v occurs in time step
218 t , and two international events data sets—GDELT [53] and ICEWS [54]—where $y_v^{(t)}$ is the number
219 of times sender-receiver pair v interacted during time step t . We obtained the matrices and random
220 masks along with the original results for both PGDS and GP-DPFA from the authors and ran the
221 PrGDS with the same MCMC settings they describe (see their paper for details [22]) and BTPF
222 (which uses variational inference) on all matrices and masks. We display the results in Fig. 4(a).

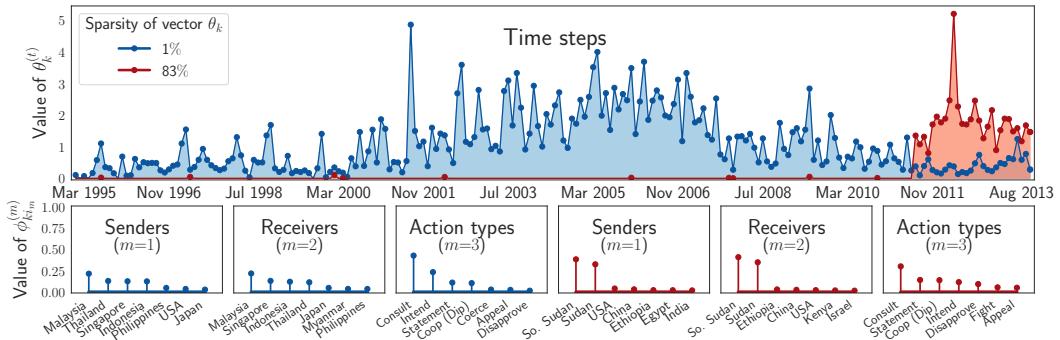
223 **Tensors.** We obtained tensor data from two international events data sets—GDELT and ICEWS—
224 wherein a count $y_{i \rightarrow j}^{(t)}$ is the number of times country i took action a towards country j during time
225 step t . These counts form a sequence of 3-mode tensors of size $V \times V \times A$ for $V = 249$ countries and
226 $A = 20$ actions. For both data sets, we treat months as time steps, where for GDELT we consider the
227 date range 2003–2008 (thus, $T = 72$) and for ICEWS we consider 1995–2013 ($T = 228$).

228 We also obtained neuroscience data of multielectrode recordings of macaque monkey motor cortices
229 from Williams et al. (2018) [55]. In this data, a count $y_{ij}^{(t)}$ is the number of times neuron i fired in trial
230 j during time step t . These counts form a sequence of matrices each of size $N \times S$ where $N = 100$ is
231 the number of neurons and $S = 1,716$ is the number of trials. We consider each time step to be a 20
232 millisecond interval which yields $T = 162$. See Fig. 4(b).

233 For each of the three tensors, we randomly generated two masks that each holdout three pairs of
234 adjacent time steps in the range $t \in [2, T-2]$ as well as the last two time steps $T-1, T$. For each
235 dynamic model, mask, and tensor, we run two independent chains of 4,000 MCMC iterations, saving
236 every 100th sample after the first 1,000 to compute heldout log-perplexity; we then also fit BPTF and
237 report the relative improvement over it for the dynamic models in Fig. 4(b).



(a) Latent activation structure inferred from macaque cortex data by the PrGDS. *Right:* Comparison of the $K \times T$ state matrix Θ inferred by the two PrGDS variants: $\epsilon_0^{(\theta)} = 0$ (sparse) vs. $\epsilon_0^{(\theta)} = 1$. The components k are sorted by which time step t had the largest $\theta_k^{(t)}$. The banded structure suggests both infer components that activate within specific short durations. White cells correspond to zeros $\theta_k^{(t)} = 0$ —the PrGDS can only represent sparse activation structure when $\epsilon_0^{(\theta)} = 0$. *Left:* Alternative visualization of the ten sparsest rows of the matrix Θ_k . Each is a component's T -length activation vector—they collectively depict localized and “bursty” neuronal spiking.



(b) Two components inferred by the sparse PrGDS from ICEWS data—the blue component was found by other models, the red component was not. The red component is specific to South Sudan, as can be seen by visualizing largest values of the factor vectors for the sender and receiver modes (*bottom row, second and third rightmost*). South Sudan was not a country until July 2011. The activation vector (*top*) is thus sparse— $\theta_k^{(t)} = 0$ at 94% of time steps (months) prior to July 2011 (83% overall). By contrast, the blue component represents Southeast Asian relations, which are persistently active. The sparse PrGDS can infer both temporally-persistent latent structure that other models infer (blue), as well as temporally-localized structure that other models do not (red).

Figure 5: Sparse representations of macaque cortex activity (Fig. 5(a)) and ICEWS events (Fig. 5(b)).

238 **Discussion.** A PrGDS variant obtained the lowest perplexity of all models in ten out of the 14 studies
 239 (i.e., subplots in Fig. 4). The sparse PrGDS obtained lower perplexity than the non-sparse PrGDS in
 240 nine studies, sometimes dramatically lower. We conjecture that the better performance of the sparse
 241 PrGDS can be explained by the expectation of future states $\theta^{(t)}$ given in Eq. (8)—when $\epsilon_0^{(\theta)} > 0$ this
 242 expectation includes an additive bias term which grows as we forecast further time steps. When $\epsilon_0^{(\theta)} = 0$
 243 the bias term vanishes and the expectation matches the analogous one for the PGDS, given in Eq. (10).

244 We also find that the sparse variant is capable of inferring a qualitatively broader range of latent
 245 structure that includes bursty and sparse latent structure. In Fig. 5(a) and Fig. 5(b) we explore some
 246 of the sparse latent structure inferred by the PrGDS from the macaque cortex data and the ICEWS
 247 event data. In the Appendix, we provide examples of how we aligned components across models and
 248 provide examples of some well-aligned ones. We found that all models inferred qualitatively similar
 249 components—however, a small number of components were unique to the sparse PrGDS. We give
 250 one such example in Fig. 5(b).

251 6 Conclusion

252 A novel modeling motif—Poisson–gamma–Poisson recursions—allows us to construct the Poisson–
 253 randomized gamma dynamical system, a tractable and flexible model family for sequentially observed
 254 count tensors. A variant of the PrGDS permits a truly sparse latent representation that is both
 255 qualitatively appealing and provides a natural inductive bias for sparse and “bursty” count time series.

256 **References**

- 257 [1] Philip A Schrot. Event data in foreign policy analysis. *Foreign Policy Analysis: Continuity*
258 *and Change in Its Second Generation*, pages 145–166, 1995.
- 259 [2] Gary King. Proper nouns and methodological propriety: Pooling dyads in international relations
260 data. *International Organization*, 55(2):497–507, 2001.
- 261 [3] Donald P Green, Soo Yeon Kim, and David H Yoon. Dirty pool. *International Organization*,
262 55(2):441–468, 2001.
- 263 [4] Paul Poast. (Mis)using dyadic data to analyze multilateral events. *Political Analysis*, 18(4):403–
264 425, 2010.
- 265 [5] Robert S Erikson, Pablo M Pinto, and Kelly T Rader. Dyadic analysis in international relations:
266 A cautionary tale. *Political Analysis*, 22(4):457–463, 2014.
- 267 [6] Brandon Stewart. Latent factor regressions for the social sciences. *Harvard University:*
268 *Department of Government Job Market Paper*, 2014.
- 269 [7] Peter D Hoff and Michael D Ward. Modeling dependencies in international relations networks.
270 *Political Analysis*, 12(2):160–175, 2004.
- 271 [8] Peter D Hoff. Multilinear tensor regression for longitudinal relational data. *The annals of*
272 *applied statistics*, 9(3):1169, 2015.
- 273 [9] A. Schein, J. Paisley, D. M. Blei, and H. Wallach. Bayesian Poisson tensor factorization for
274 inferring multilateral relations from sparse dyadic event counts. In *Proceedings of the Twenty-*
275 *First ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages
276 1045–1054, 2015.
- 277 [10] Peter D Hoff et al. Equivariant and scale-free Tucker decomposition models. *Bayesian Analysis*,
278 11(3):627–648, 2016.
- 279 [11] Aaron Schein, Mingyuan Zhou, David M. Blei, and Hanna Wallach. Bayesian Poisson Tucker
280 decomposition for learning the structure of international relations. In *Proceedings of the 33rd*
281 *International Conference on Machine Learning*, 2016.
- 282 [12] Jon Kleinberg. Bursty and hierarchical structure in streams. *Data Mining and Knowledge*
283 *Discovery*, 7(4):373–397, 2003.
- 284 [13] E. C. Chi and T. G. Kolda. On tensors, sparsity, and nonnegative factorizations. *SIAM Journal*
285 *on Matrix Analysis and Applications*, 33(4):1272–1299, 2012.
- 286 [14] Tsuyoshi Kunihama and David B Dunson. Bayesian modeling of temporal dependence in large
287 sparse contingency tables. *Journal of the American Statistical Association*, 108(504):1324–1338,
288 2013.
- 289 [15] J. Canny. GaP: A factor model for discrete data. In *Proceedings of the 27th Annual International*
290 *ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 122–129,
291 2004.
- 292 [16] D. B. Dunson and A. H. Herring. Bayesian latent variable models for mixed discrete outcomes.
293 *Biostatistics*, 6(1):11–25, 2005.
- 294 [17] Michalis K. Titsias. The infinite gamma–Poisson feature model. In *Advances in Neural*
295 *Information Processing Systems 21*, pages 1513–1520, 2008.
- 296 [18] A. T. Cemgil. Bayesian inference for nonnegative matrix factorisation models. *Computational*
297 *Intelligence and Neuroscience*, 2009.
- 298 [19] M. Zhou, L. Hannah, D. Dunson, and L. Carin. Beta-negative binomial process and Poisson
299 factor analysis. In *Proceedings of the 15th International Conference on Artificial Intelligence*
300 *and Statistics*, 2012.

- 301 [20] Prem K Gopalan and David M Blei. Efficient discovery of overlapping communities in massive
302 networks. *Proceedings of the National Academy of Sciences*, 110(36):14534–14539, 2013.
- 303 [21] Beyza Ermis and A Taylan Cemgil. A bayesian tensor factorization model via variational
304 inference for link prediction. *arXiv preprint arXiv:1409.8276*, 2014.
- 305 [22] Aaron Schein, Hanna Wallach, and Mingyuan Zhou. Poisson-gamma dynamical systems. In
306 *Advances in Neural Information Processing Systems*, pages 5005–5013, 2016.
- 307 [23] Lin Yuan and John D Kalbfleisch. On the Bessel distribution and related problems. *Annals of*
308 *the Institute of Statistical Mathematics*, 52(3):438–447, 2000.
- 309 [24] R. Harshman. Foundations of the PARAFAC procedure: Models and conditions for an “ex-
310 planatory” multimodal factor analysis. *UCLA Working Papers in Phonetics*, 16:1–84, 1970.
- 311 [25] Roman N Makarov and Devin Glew. Exact simulation of Bessel diffusions. *Monte Carlo*
312 *Methods and Applications*, 16(3-4):283–306, 2010.
- 313 [26] Milton Abramowitz and Irene A Stegun. *Handbook of mathematical functions: with formulas,*
314 *graphs, and mathematical tables*, volume 55. Courier Corporation, 1965.
- 315 [27] Mingyuan Zhou, Yulai Cong, and Bo Chen. Augmentable gamma belief networks. *The Journal*
316 *of Machine Learning Research*, 17(1):5656–5699, 2016.
- 317 [28] Chengyue Gong et al. Deep dynamic Poisson factorization model. In *Advances in Neural*
318 *Information Processing Systems*, pages 1666–1674, 2017.
- 319 [29] Dandan Guo, Bo Chen, Hao Zhang, and Mingyuan Zhou. Deep Poisson gamma dynamical
320 systems. In *Advances in Neural Information Processing Systems*, pages 8442–8452, 2018.
- 321 [30] A. Acharya, J. Ghosh, and M. Zhou. Nonparametric Bayesian factor analysis for dynamic
322 count matrices. *Proceedings of the 18th International Conference on Artificial Intelligence and*
323 *Statistics*, 2015.
- 324 [31] Sikun Yang and Heinz Koepll. Dependent relational gamma process models for longitudinal
325 networks. In *International Conference on Machine Learning*, pages 5547–5556, 2018.
- 326 [32] M. Zhou and L. Carin. Augment-and-conquer negative binomial processes. In *Advances in*
327 *Neural Information Processing Systems Twenty-Five*, pages 2546–2554, 2012.
- 328 [33] Mingyuan Zhou, Yulai Cong, and Bo Chen. The Poisson gamma belief network. In *Advances*
329 *in Neural Information Processing Systems*, pages 3043–3051, 2015.
- 330 [34] A Taylan Cemgil and Onur Dikmen. Conjugate gamma Markov random fields for modelling
331 nonstationary sources. In *International Conference on Independent Component Analysis and*
332 *Signal Separation*, pages 697–705. Springer, 2007.
- 333 [35] Cédric Févotte, Jonathan Le Roux, and John R Hershey. Non-negative dynamical system with
334 application to speech and audio. In *2013 IEEE International Conference on Acoustics, Speech*
335 *and Signal Processing*, pages 3158–3162. IEEE, 2013.
- 336 [36] Ghassen Jerfel, Mehmet E Basbug, and Barbara E Engelhardt. Dynamic collaborative filtering
337 with compound Poisson factorization. *arXiv preprint arXiv:1608.04839*, 2016.
- 338 [37] Rajesh Ranganath, Linpeng Tang, Laurent Charlin, and David Blei. Deep exponential families.
339 In *Artificial Intelligence and Statistics*, pages 762–771, 2015.
- 340 [38] Robert GD Steel. Relation between Poisson and multinomial distributions. 1953.
- 341 [39] Anirban Bhattacharya and David B Dunson. Simplex factor models for multivariate unordered
342 categorical data. *Journal of the American Statistical Association*, 107(497):362–377, 2012.
- 343 [40] Jakob H Macke, Lars Buesing, John P Cunningham, M Yu Byron, Krishna V Shenoy, and
344 Maneesh Sahani. Empirical models of spiking in neural populations. In *Advances in neural*
345 *information processing systems*, pages 1350–1358, 2011.

- 346 [41] Rudolph E Kalman and Richard S Bucy. New results in linear filtering and prediction theory.
347 *Journal of basic engineering*, 83(1):95–108, 1961.
- 348 [42] Zoubin Ghahramani and Sam T Roweis. Learning nonlinear dynamical systems using an EM
349 algorithm. In *Advances in neural information processing systems*, pages 431–437, 1999.
- 350 [43] John Ashworth Nelder and Robert WM Wedderburn. Generalized linear models. *Journal of the
351 Royal Statistical Society: Series A (General)*, 135(3):370–384, 1972.
- 352 [44] Laurent Charlin, Rajesh Ranganath, James McInerney, and David M Blei. Dynamic Poisson
353 factorization. In *Proceedings of the 9th ACM Conference on Recommender Systems*, pages
354 155–162. ACM, 2015.
- 355 [45] Patrick T Brandt and Todd Sandler. A Bayesian Poisson vector autoregression model. *Political
356 Analysis*, 20(3):292–315, 2012.
- 357 [46] Charles Blundell, Jeff Beck, and Katherine A Heller. Modelling reciprocating relationships with
358 Hawkes processes. In *Advances in Neural Information Processing Systems*, pages 2600–2608,
359 2012.
- 360 [47] Aleksandr Simma and Michael I Jordan. Modeling events with cascades of Poisson processes.
361 *arXiv preprint arXiv:1203.3516*, 2012.
- 362 [48] Scott Linderman and Ryan Adams. Discovering latent network structure in point process data.
363 In *International Conference on Machine Learning*, pages 1413–1421, 2014.
- 364 [49] Luc Devroye. Simulating Bessel random variables. *Statistics & Probability Letters*, 57(3):249–
365 257, 2002.
- 366 [50] NeurIPS corpus. UCI Machine Learning Repository.
- 367 [51] dblp computer science bibliography. <http://dblp.uni-trier.de/>.
- 368 [52] State of the Union Addresses (1790-2006) by United States Presidents. [https://www.gutenberg.org/ebooks/5050?msg=wELCOME_STRANGER](https://www.gutenberg.org/ebooks/5050?msg=welcome_stranger).
- 370 [53] K. Leetaru and P. Schrodt. GDELT: Global data on events, location, and tone, 1979–2012.
371 Working paper, 2013.
- 372 [54] E. Boschee, J. Lautenschlager, S. O'Brien, S. Shellman, J. Starz, and M. Ward. ICEWS coded
373 event data. Harvard Dataverse, 2015. V10.
- 374 [55] Alex H Williams, Tony Hyun Kim, Forea Wang, Saurabh Vyas, Stephen I Ryu, Krishna V
375 Shenoy, Mark Schnitzer, Tamara G Kolda, and Surya Ganguli. Unsupervised discovery of
376 demixed, low-dimensional neural dynamics across multiple timescales through tensor compo-
377 nent analysis. *Neuron*, 98(6):1099–1115, 2018.