

# Dynamic Poisson Factor Analysis

**Abstract**—We introduce a novel dynamic model for discrete time-series data, in which the temporal sampling may be nonuniform. The model is specified by constructing a hierarchy of Poisson factor analysis blocks, one for the *transitions* between latent states and the other for the *emissions* between latent states and observations. Latent variables are binary and linked to Poisson factor analysis via Bernoulli-Poisson specifications. The model is derived for count data but can be readily modified for binary observations. We derive efficient inference via Markov chain Monte Carlo, that scales with the number of non-zeros in the data and latent binary states, yielding significant acceleration compared to related models. Experimental results on benchmark data show the proposed model achieves state-of-the-art predictive performance. Additional experiments on microbiome data demonstrate applicability of the proposed model to interesting problems in computational biology where interpretability is of utmost importance.

## I. INTRODUCTION

Probabilistic models for high-dimensional time series have long been an area of significant interest in machine learning. The applicability of these models spans data from different domains, such as music sequences, text streams and time-series of molecular data in computational biology. Among this prior work, Hidden Markov Models (HMMs) [1] and the Linear Dynamical System (LDS) [2] are particularly well understood. However, in some cases these models are limited by the type of dynamic structures they can capture. In fact, real-world time-series data are usually described by complex nonlinear temporal dependencies, while traditional LDS models are restricted to latent representations described by linear dynamics. Models with discrete latent spaces, such as the HMM, are often specified as mixture models that represent the history of a time-series using multinomial distributions.

A newer class of time-series models, better suited to model complex (nonlinear) probability distributions over high-dimensional sequences, rely either on Recurrent Neural Networks (RNNs) [3]–[6], Restricted Boltzmann Machines (RBMs) [7]–[11] or Sigmoid Belief Networks (SBNs) [12]. The Temporal Restricted Boltzmann Machine (TRBM) [8] and the Temporal Sigmoid Belief Network (TSBN) [12], for instance, consist of a sequence of RBMs or SBNs, respectively, where the state of the current RBM or SBN is stochastically determined by previous RBMs or SBNs. Inference approaches for these models are non-trivial. Approximate procedures based mainly on Variational Bayes principles have been proposed and scale well to large datasets [8], [11], [12].

In the context of count data, multi-layer directed models are becoming increasingly popular [12]–[15]. In these models, the likelihood function connecting latent variables to observations is often specified in terms of Poisson distributions or softmax link functions, whereas latent (often deep) layers of the

model are described by binary variables, capturing nonlinear dependencies and often modeled via SBNs [13], [14]. For time-series data, a similar idea has been leveraged, where a discrete transition model akin to HMMs connects current binary latent variables to previous ones (at earlier time points), but deviates from HMMs by specifying transitions via SBNs [12], not multinomial distributions.

In the work presented here, we propose a model for time-series data where both emission and transition models are specified in terms of Poisson distributions, borrowing ideas from deep Poisson factor models [15]. In particular, the Poisson-based transition model treats transition *strengths* as latent counts that are transformed to binary variables via the Bernoulli-Poisson Link (BPL) [16], a recently proposed alternative to the sigmoid link function. The BPL yields efficient learning and inference algorithms [15], [16]. In particular, the key advantage of modeling transitions with the BPL is that learning and inference scales with the number of non-zero latent variables, as opposed to the number of states (like is the case of TRBM or TSBN models), where the sigmoid link function is employed [8], [12].

The main contributions of this work are:

- (i) We develop a dynamic model for times-series data with count observations based on Poisson factor analysis. The Bernoulli-Poisson link formulation can be readily accommodated to time-series with binary observations.
- (ii) Unlike most previously proposed models, our formulation easily allows for modeling data *nonuniformly* sampled in time.
- (iii) An efficient sampling inference procedure is developed, scaling with the number of non-zeros in the data and binary latent variables, where latent counts and binary variables in the model can be updated in block. This allows our implementation to benefit from significant parallelization using GPUs (demonstrated in our experiments).
- (iv) Results on benchmark and real data highlight the benefits of our modeling strategy from the stand-points of performance and interpretability; this is demonstrated via analysis of a new dynamic microbiome dataset.

## II. DYNAMIC POISSON FACTOR ANALYSIS

Assume observed counts for  $N$  time-series, where the vector of counts at each time point is of dimension  $M$ . The vector of counts at time  $t$  for the  $n$ th time series is denoted as  $\mathbf{x}_{nt} \in \mathbb{N}^M$ . Akin to HMMs, we model the dynamics of the time-series by imposing a *transition* model on latent variables,

$\mathbf{h}_{nt} \in \{0, 1\}^K$ , where the state of the current latent variable (at time  $t$ ) depends on the previous state (at time  $t-1$ ), and  $K$  is the number of latent variables. The model allows  $2^K$  different realizations of  $\mathbf{h}_{nt}$ , yielding  $2^K$  states. However, the proposed model characterizes transition dynamics by more than just these discrete states (moving beyond an HMM), yielding improved modeling flexibility and results, as detailed below.

The joint probability of the  $n$ th observation at time  $t$  is

$$p(\mathbf{X}_n, \mathbf{H}_n | \Xi, \Omega) = p(\mathbf{h}_{n0}) \prod_{t=1}^{T_n} p_{\Xi}(\mathbf{x}_{nt} | \mathbf{h}_{nt}) p_{\Omega}(\mathbf{h}_{nt} | \mathbf{h}_{nt-1}), \quad (1)$$

where  $\mathbf{X}_n = [\mathbf{x}_{n1}, \dots, \mathbf{x}_{nT_n}]$ ,  $\mathbf{H}_n = [\mathbf{h}_{n1}, \dots, \mathbf{h}_{nT_n}]$  and  $T_n$  is the length of the  $n$ th time-series. Further,  $p(\mathbf{h}_{n0})$  is the prior for the initial value of the latent variables,  $p_{\Xi}(\mathbf{x}_{nt} | \mathbf{h}_{nt})$  is the *emission* model with parameters  $\Xi$  and  $p_{\Omega}(\mathbf{h}_{nt} | \mathbf{h}_{nt-1})$  is the transition model with parameters  $\Omega$ .

#### A. Emission model

For the observed vector  $\mathbf{x}_{nt}$ , containing counts of  $M$  entities (e.g., a vocabulary of  $M$  words), we impose the following *emission* model

$$\mathbf{x}_{nt} \sim \text{Poisson}(\Psi(\theta_{nt} \circ \mathbf{h}_{nt})), \quad (2)$$

where  $\Psi \in \mathbb{R}_+^{M \times K}$  is the *global* factor loadings matrix with  $K$  factors, shared by all time-series and time points;  $\theta_{nt} \in \mathbb{R}_+^K$  and  $\mathbf{h}_{nt} \in \{0, 1\}^K$  are *local* variables representing factor intensities and activations, respectively. Note that entries of  $\mathbf{h}_{nt}$  indicate which factors are active for observation  $n$  at time  $t$ , i.e., they define the *state* to which  $\mathbf{x}_{nt}$  belongs. Symbol  $\circ$  represents the element-wise (Hadamard) product.

Note that the Poisson emission parameters for data  $n$  are not just dependent on the state at time  $t$ ,  $\mathbf{h}_{nt}$ . Binary activations,  $\mathbf{h}_{nt}$ , impose which columns of  $\Psi$  are employed to represent the Poisson rate, and the nonnegative intensities,  $\theta_{nt}$ , provide temporal- and data-dependent scaling; in our experiments we have found this modeling flexibility to be important, particularly on real data, e.g., the motivating microbiome data.

The model in (2) may be rewritten as

$$\begin{aligned} x_{mnt} &= \sum_{k=1}^K x_{mknt}, \\ x_{mknt} &\sim \text{Poisson}(\lambda_{mknt}), \\ \lambda_{mknt} &= \psi_{mk} \theta_{knt} h_{knt}, \end{aligned} \quad (3)$$

where  $x_{mnt}$  is component  $m$  of vector  $\mathbf{x}_{nt}$ ,  $\psi_{mk}$  is component  $m$  of  $\psi_k$ ,  $\psi_k$  is column  $k$  of  $\Psi$ ,  $\theta_{knt}$  is component  $k$  of vector  $\theta_n$ , and  $h_{knt}$  is component  $k$  of vector  $\mathbf{h}_n$ . In (3) we have used the additive property of the Poisson distribution to decompose the  $m$ th observed count of  $\mathbf{x}_{nt}$  as  $K$  latent counts,  $\{x_{mknt}\}_{k=1}^K$ . This decomposition allows derivation of efficient inference for the entire model, as discussed in Section III.

We specify prior distributions for the model in (2) as previously described [17], i.e.,

$$\begin{aligned} \psi_k &\sim \text{Dirichlet}(\eta_{\psi} \mathbf{1}_M), \\ \theta_{knt} &\sim \text{Gamma}(r_k, b_{\theta}), \\ h_{knt} &\sim \text{Bernoulli}(\pi_{knt}), \end{aligned} \quad (4)$$

where  $\mathbf{1}_M$  is an  $M$ -dimensional vector of all-ones. Favoring simplicity, we let  $\eta_{\psi} = 1/K$ ,  $b_{\theta} = 0.5$  and  $r_k \sim \text{Gamma}(1, 1)$ . Prior distributions for  $\eta_{\psi}$  and  $b_{\theta}$  that result in closed form conditionals exist, and can be used if desired; see for instance [18] for  $\eta_{\psi}$ , and [19] for  $b_{\theta}$ .

The hierarchical model implied by (2) and (4), known as Poisson Factor Analysis (PFA), corresponds to the emission model,  $p_{\Xi}(\mathbf{x}_{nt} | \mathbf{h}_{nt})$ , succinctly expressed in (1), with parameters  $\Xi = \{\Psi, \theta_{nt}, r_k\}$ , for  $n = 1, \dots, N$ ,  $t = 1, \dots, T_n$  and  $k = 1, \dots, K$ . These parameters can be interpreted in the context of topic modeling as follows:  $\Psi$  is a loadings matrix whose columns encode  $K$  topics (distributions over  $M$  words), that capture the correlation structure of observed variables; binary latent activations  $\mathbf{h}_{nt}$  select which topics are used in time-series  $n$  at time  $t$ ; and  $\theta_{nt}$ , encode the intensities with which each topic is manifested in observation  $\mathbf{x}_{nt}$ . Interestingly, the PFA is closely related to other well-known topic modeling approaches, such as latent Dirichlet allocation, hierarchical Dirichlet processes and focused topic models [19].

#### B. Transition model

The most unique aspect of the proposed model is how state transitions are modeled, and how this allows nonuniform temporal sampling. In order to specify our model for latent variable transitions with respect to time, we first introduce the Bernoulli-Poisson link [16], a recently proposed probabilistic link function particularly useful at relating binary and count variables. Specifically, for a binary vector  $\mathbf{h}_{nt}$  with elements  $h_{knt}$ ,

$$\begin{aligned} h_{knt} &= 1(z_{knt} > 0), \\ z_{knt} &\sim \text{Poisson}(\tilde{\lambda}_{knt}), \end{aligned} \quad (5)$$

where  $z_{knt}$  is a latent count associated with binary variable  $h_{knt}$ , parameterized by a Poisson distribution with rate  $\tilde{\lambda}_{knt}$ . The function  $1(\cdot)$  is defined as  $1(\cdot) = 1$  if the argument holds, and  $1(\cdot) = 0$ , otherwise. The model in (5), denoted here for short as  $\mathbf{h}_{nt} \sim \text{BPL}(\tilde{\lambda}_{nt})$ , for  $\tilde{\lambda}_{nt} \in \mathbb{R}_+^K$  with elements  $\tilde{\lambda}_{knt}$ , has the interesting property that

$$\begin{aligned} p(h_{knt} = 1) &= \text{Bernoulli}(\pi_{knt}), \\ \pi_{knt} &= 1 - \exp(-\tilde{\lambda}_{knt}). \end{aligned} \quad (6)$$

The result in (6) can be shown by marginalizing out latent counts,  $z_{knt}$ , in (5). In fact, to sample  $h_{knt}$  we do not need to instantiate latent count  $z_{knt}$ , but the rate of its underlying Poisson distribution,  $\tilde{\lambda}_{knt}$ .

The distribution implied by (6) is reminiscent of the complementary log-log link function [20], [21], where  $\tilde{\lambda} = \exp(-u)$

and  $u \in \mathbb{R}$ . The logistic link function used in RBMs and SBNs is symmetric around the origin,  $u = 0$ , with  $p(h = 1) = \text{Bernoulli}(\pi)$ , where  $\pi = 1/(1 + \exp(-u))$ ; by contrast, the proposed BPL link is *asymmetric*, which makes it appropriate for very sparse settings, where the proportion of zeros is large. In our case, this is particularly useful, because we can increase the number of latent variables without forcing the model to increase the number of states *a priori*.

Having defined the Bernoulli-Poisson link, it becomes clear how to specify a transition model using the same Poisson factor analysis framework used for the emission model. Specifically, we write

$$\begin{aligned} \mathbf{h}_{nt} &\sim \text{BPL} \left( \tau_{nt}^{-1} \Phi(\mathbf{w}_{nt-1} \circ \mathbf{h}_{nt-1}) + \tilde{\lambda}_0 \right), \\ \phi_k &\sim \text{Dirichlet}(\eta_\phi \mathbf{1}_K), \\ w_{knt-1} &\sim \text{Gamma}(s_k, b_w), \end{aligned} \quad (7)$$

where  $\tilde{\lambda}_{nt} = \tau_{nt}^{-1} \Phi(\mathbf{w}_{nt-1} \circ \mathbf{h}_{nt-1}) + \tilde{\lambda}_0$  as in (5) and (6),  $\phi_k$  is a column of  $\Phi$  (transition factor matrix),  $w_{knt-1}$  is an element of  $\mathbf{w}_{nt-1}$  (*local* variable representing transition factor intensity), and  $\eta_\phi$ ,  $b_w$  and  $s_k$  are specified in a similar fashion to the emission model in (4). Bias term,  $\tilde{\lambda}_0$ , controls the base rate of the Poisson distribution in (5) and is specified below. Parameter  $\tau_{nt}$  is the time difference between observations  $t$  and  $t-1$ , for time-series  $n$ ;  $\tau_{nt}$  can vary with  $n$  and  $t$ , allowing *nonuniform* temporal sampling.

The specification in (7) corresponds to the transition model,  $p_\Omega(\mathbf{h}_{nt}|\mathbf{h}_{nt-1})$ , in (1), with parameters  $\Omega = \{\Phi, \mathbf{w}_{nt}, s_k, \lambda_0\}$ , for  $n = 1, \dots, N$ ,  $t = 1, \dots, T_n$  and  $k = 1, \dots, K$ . These parameters have clear interpretations. For instance,  $\Phi$  is a transition matrix whose columns,  $\phi_k$ , encode distinct transition *templates*. These templates are  $K$ -dimensional probability vectors and can be viewed as distributions over latent binary variable activations,  $\mathbf{h}_{nt}$ . Each template defines a particular activation pattern; some latent variables are co-active with high probability, while others are jointly absent, *i.e.*, they define correlation structure among elements of  $\mathbf{h}_{nt}$ . Interestingly, templates at time  $t$  are selected by previous activations  $\mathbf{h}_{nt-1}$ , and regulated (weighted) by previous intensities  $\mathbf{w}_{nt-1}$ . This implies that at time  $t$  the transition statistics are not only dependent on the  $2^K$  discrete states  $\mathbf{h}_{nt}$ , but scale these states with weights  $\mathbf{w}_{nt-1}$ , adding modeling flexibility analogous to that in the emission model discussed above. In addition, the correlation between two adjacent time-points decays inversely with proximity in time, *i.e.*, as  $\tau_{nt}$  increases, the dependency between latent variables  $\mathbf{h}_{nt}$  and  $\mathbf{h}_{nt-1}$  decreases. When  $\tau_{nt}$  is large enough, binary activations loose their time dependency and effectively become a stochastic function of  $\tilde{\lambda}_0$ , in fact, from (6)

$$p(h_{knt} = 1) = \text{Bernoulli} \left( 1 - \exp(-\tilde{\lambda}_{k0}) \right),$$

where  $\tilde{\lambda}_{k0}$  is an element of  $\tilde{\lambda}_0$ .

Note that the fact that intensities are time-dependent adds flexibility to the model compared to a simplified specification

where these intensities are global parameters, *i.e.*,  $w_{kn}$  instead of  $w_{knt}$ . We observed empirically that the simplified model (with time-independent weights) does not perform as well as the specification in (7), however it is worth mentioning that time-dependent intensities may undermine the contribution of  $\tau_{nt}$  to the transition model.

The emission model is completed by specifying a prior distribution for the initial state of the latent variables of the dynamic model,  $p(\mathbf{h}_{n0})$ , in (1). We let

$$\begin{aligned} h_{kn0} &\sim \text{BPL} \left( \tilde{\lambda}_{k0} \right), \\ \tilde{\lambda}_{k0} &\sim \text{Gamma}(a, b). \end{aligned} \quad (8)$$

For simplicity, we let  $a = b = 1$ , so that elements of  $\mathbf{h}_{n0}$  are approximately uniform. Note that (8) is a special case of (7) because since we do not have past information about  $\mathbf{h}_{n0}$ , conceptually,  $\tau_{n0} \rightarrow \infty$ .

In summary, the joint distribution for our dynamic Poisson factor model in (1) is fully specified by the emission model implied by (2) and (4), the transition model implied by (5) and (7), and the initial-state probability as in (8). Figure 1 shows a graphical representation of the model in terms of emission and transition rates,  $\lambda_{nt} = \Psi(\theta_{nt} \circ \mathbf{h}_{nt})$  and  $\tilde{\lambda}_{nt} = \tau_{nt}^{-1} \Phi(\mathbf{w}_{nt-1} \circ \mathbf{h}_{nt-1}) + \tilde{\lambda}_0$ , respectively.

### C. Binary data

We can easily extend the dynamic Poisson factor model described above to model binary time-series data, by leveraging the same type of construction used for the transition model, *i.e.*, for  $\mathbf{x}_{nt} \in \{0, 1\}^M$ , we can let

$$\mathbf{x}_{nt} \sim \text{BPL}(\Psi(\theta_{nt} \circ \mathbf{h}_{nt})),$$

as in (7) but with prior distributions defined as in (4).

## III. LEARNING AND INFERENCE

The dynamic Poisson factor model defined in the previous section has the convenient property of having all its conditional posteriors available in closed form, due to local conjugacy. In this paper, we focus on Markov Chain Monte Carlo (MCMC) via Gibbs sampling for learning and inference. Stochastic variational inference may be also readily implemented using ideas from [15]. Other alternatives for scaling up inference based on gradient-based approaches [22]–[24] are also amenable to our model specification, however beyond the scope of this paper.

During the *learning* phase, Gibbs sampling for the model in (2), (4), (5), (7) and (8) involves sampling in sequence from the conditional posterior of all the global parameters of the model, namely,  $\{\Psi, r_k, \Phi, s_k, \tilde{\lambda}_{k0}\}$ , for  $k = 1, \dots, K$ . During the *inference* phase, we sample from the conditional posterior of all local parameters, while also sampling from the *fixed* conditional posterior of the global parameters given the training data (obtained during learning). The local parameters include all latent variable activations,  $\{\mathbf{h}_{nt}\}$ , intensities for the emission model,  $\theta_{nt}$ , and intensities for the transition model,  $\mathbf{w}_{nt}$ , where  $n = 1, \dots, N$ ,  $t = 1, \dots, T_n$ . For prediction tasks, we perform learning on a training set, then perform inference

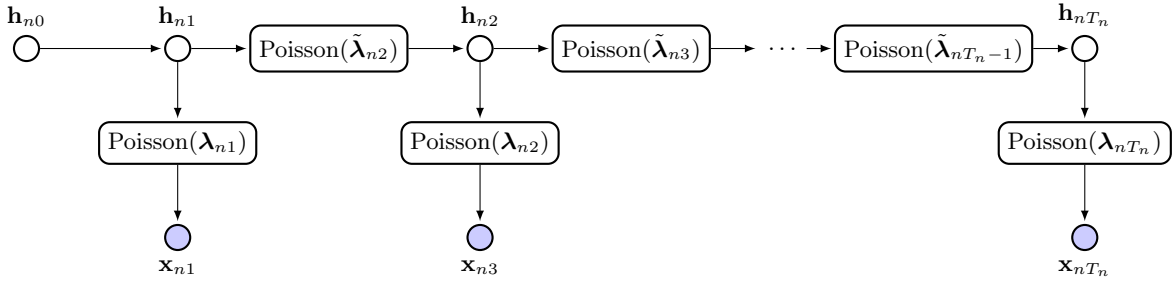


Fig. 1. Graphical model for dynamic Poisson factor analysis. For the transitions,  $\lambda_{nt} = \Psi(\theta_{nt} \circ \mathbf{h}_{nt})$ , and for the transitions,  $\tilde{\lambda}_{nt} = \tau_{nt}^{-1} \Phi(\mathbf{w}_{nt-1} \circ \mathbf{h}_{nt-1}) + \tilde{\lambda}_0$ , as defined in (2) and (7), respectively. Filled and empty nodes represent observed and latent variables, respectively.

on the test set, while sampling from the learned conditional posterior of the global parameters conditioned only on training data. In this way, we can benefit from model averaging at test time.

The hyperparameters of the model are set to fixed values:  $\eta_\psi = \eta_\phi = 1/K$ ,  $b_\theta = b_w = 0.5$  and  $a_0 = b_0 = 1$ . Note that priors for  $\eta$ ,  $b$ ,  $a_0$  and  $b_0$  exist that result in Gibbs-style updates, and can be readily incorporated into the model if desired; however, our priority is keeping the model simple without sacrificing flexibility. The most unique conditional posteriors are shown below

#### Latent counts:

$$x_{mknt} \sim \text{Multinomial}\left(x_{mnt}, [\hat{\lambda}_{m1nt}, \dots, \hat{\lambda}_{mKnt}]\right), \quad (9)$$

where  $\hat{\lambda}_{m \cdot nt} = \lambda_{m \cdot nt} / \sum_k \lambda_{mknt}$ ,  $\lambda_{mknt} = \psi_{mk} \theta_{knt} h_{knt}$ .

#### Factor loadings columns:

$$\psi_k \sim \text{Dirichlet}(\eta_\psi + x_{1k\cdot}, \dots, \eta_\psi + x_{Mk\cdot}),$$

where  $x_{mk\cdot} = \sum_{n,t} x_{mknt}$ .

#### Factor intensities:

$$\theta_{knt} \sim \text{Gamma}(r_k h_{knt} + x_{\cdot knt}, b_\theta), \quad (10)$$

$$r_k \sim \text{Gamma}\left(1 + \sum_{n,t} \ell_{knt}, 1 - \sum_{n,t} h_{knt} \log(1 - b_\theta)\right),$$

where  $x_{\cdot knt} = \sum_m x_{mknt}$  and  $\ell_{knt} \sim \text{CRT}(x_{\cdot ntk}, r_k)$  is the Chinese Restaurant Table (CRT) distribution [19]. Note that the conditional update for  $r_k$  is obtained by leveraging the gamma-Poisson conjugacy and the data augmentation scheme described in [19]. Briefly, the gamma-Poisson construction of  $x \sim \text{Poisson}(\lambda)$ ,  $\lambda \sim \text{Gamma}(r, b/(1-b))$  can be represented as  $m = \sum_{t=1}^L \ell_t$ ,  $\ell_t \sim \text{Log}(b)$ ,  $L \sim \text{Poisson}(-r \ln(1-b))$ , where  $\text{Log}(\cdot)$  denote the logarithmic distribution [25].

#### Factor activations:

$$h_{knt} \sim \delta(x_{\cdot knt} = 0 \wedge z_{kt} = 0) \text{Bernoulli}(\hat{\pi}_{knt}) + \delta(x_{\cdot knt} > 0 \vee z_{kt} > 0),$$

where

$$\hat{\pi}_{knt} = \frac{\tilde{\pi}_{knt}}{\tilde{\pi}_{knt} + (1 - \pi_{knt})},$$

$$\tilde{\pi}_{knt} = \pi_{knt} (1 - b_\theta)^{r_k} (1 - b_w)^{s_k}.$$

#### Factor intensities for transition model:

$$w_{knt} \sim \text{Gamma}\left(s_k h_{knt} + z_{knt}, \tilde{b}_w(\tau_{nt})\right),$$

$$s_k \sim \text{Gamma}\left(1 + \sum_{n,t} \ell'_{knt}, 1 - \sum_{n,t} h_{knt} \log(1 - \tilde{b}_w(\tau_{nt}))\right),$$

$$\tilde{b}_w(\tau_{nt}) = \frac{b_w}{b_w + (1 - b_w) \tau_{nt}}$$

Similar to (10),  $\ell'_{knt} \sim \text{CRT}(z_{knt}, r_k)$  is drawn from the CRT distribution.

#### Latent counts for transition model:

$$z_{knt} \sim \delta(h_{knt} = 1) \text{Poisson}_+(\tilde{\lambda}_{knt}),$$

where  $\text{Poisson}_+(\tilde{\lambda})$  is the truncated Poisson distribution with rate  $\tilde{\lambda}$ . Note that from (5), if  $h_{knt} = 1$  when  $z_{knt} > 0$ , then conditioned on  $h_{knt} = 1$ , latent count variable  $z_{knt} | h_{knt} = 1$  is truncated Poisson with rate  $\tilde{\lambda}_{knt}$ .

In all the expressions above, summations over  $n$ ,  $t$  and  $m$  run with  $n = 1, \dots, N$ ,  $t = 1, \dots, T_n$  and  $m = 1, \dots, M$ , respectively.

The conditional posteriors for  $\Phi$ ,  $\pi_0$  are similar to the ones presented above, thus omitted here for conciseness. In practice, we initialize model parameters at random from their corresponding prior distributions.

An important property from our dynamic Poisson factor model is that inference does not scale with the size of the data,  $\{M, N, T_n\}$ , for  $n = 1, \dots, N$ , and the number of factors,  $K$ , but as a function of their non-zero elements, which is tremendously advantageous in cases where the data is sparse, which is often the case. In order to show that this scaling behavior holds, it is enough to see that by construction, from (3), if  $x_{mnt} = \sum_{k=1}^K x_{mknt} = 0$  (or  $z_{mnt}$ ), thus  $x_{mkn} = 0, \forall k$  with probability 1. From (5) we see that if  $h_{knt} = 0$  then  $z_{knt} = 0$  with probability 1. As a result,

update equations for all parameters of the model except for binary activations  $\mathbf{h}_{nt}$ , depend only on non-zero elements of  $\mathbf{x}_{nt}$  and  $\mathbf{z}_{nt}$ . Updates for the binary variables can be cheaply obtained in block from  $h_{knt} \sim \text{Bernoulli}(\pi_{knt})$  via  $\lambda_{knt}$ , as previously described in (6).

The most computationally expensive operation in our inference procedure is sampling from the multinomial distribution in (9), to obtain latent counts for  $x_{mknt}$  (and  $z_{knt}$ ),  $\forall m, k, n, t$ . Fortunately, conditioned on data  $\mathbf{x}_{nt}$  and emission rates  $\boldsymbol{\lambda}_{nt}$  (and  $\tilde{\boldsymbol{\lambda}}_{nt}$  for transitions),  $\forall n, t$ , these can be sampled in block using a heavily parallelized implementation of the multinomial sampler via GPUs. Results of this efficient implementation are shown in the experiments.

#### IV. RELATED WORK

The RBM has been widely used as a building block to learn the sequential dependencies in time-series data, *e.g.*, the conditional-RBM-related models [7], [26], [27], and the temporal RBM [8]. Exact inference in these partially directed models is possible via the recurrent temporal RBM [9], and further extended to learn the dependency structure within observations [11]. The model presented here is fully directed, allows for exact inference as all conditional posterior distributions are available in closed-form, and learns the dependency within observations by imposing correlation structure via factor loadings matrix,  $\Psi$ .

Our work focuses on directed generative models. Most of the existing work on directed models is based on the Sigmoid Belief Network (SBN) [28], which only recently has shown potential for building large multi-layer (deep) and dynamic models [12], [13], [29], [30]. Our model is most related with the Temporal SBN (TSBN) of [12], in which the emission model is an SBN or a softmax belief network for binary and count time-series, respectively, and the transition model is an SBN. The TSBN delivers fast inference via the Neural Variational Inference and Learning (NVIL) algorithm [13]. Our model is different from TSBN in three ways: (i) We use Bernoulli-Poisson links, not sigmoid links, for fast inference. (ii) We model observed counts directly using Poisson distributions, as opposed to observed count proportions like in TSBN, via softmax link. (iii) Our inference approach scales with the number of non-zeros in the data and binary latent variables, which is not possible for models based on SBNs (or RBMs).

The work presented here is closely related to the deep Poisson factor model [15], a multi-layer (deep) topic model in which adjacent layers are connected via BPLs similar to our model. The main differentiator between the deep PFA and our dynamic PFA is that in the former, BPLs are introduced as a way to model correlation across latent variables whereas in the latter (ours), BPLs provide us with a way to model correlation across observations, by coupling adjacent time-points (transition model). Our approach builds upon this framework by extending the functionality of Poisson factor models in general, to model time-series data. We could in principle specify a *deep* version of our model by letting emission and/or transition models have multi-layer structures as in deep PFA,

for improved flexibility and representation ability. We leave this interesting direction as future work.

To the best of our knowledge there is only one existing approach for time-series modeling based on Poisson factor analysis [31]. In their work, dynamics are imposed using a linear model on the intensities,  $\theta_{nt}$ , in terms on previous intensities  $\theta_{nt-1}$  via a gamma distribution specification,  $w_{knt} \sim \text{Gamma}(w_{knt-1}, b_\theta)$ . Our model is different in that we learn correlations across binary latent variables using  $\Phi$ , whereas in [31] latent variables are *i.i.d. a priori* and restricted to linear dynamics. Unlike ours their model does not model latent variable activations, *i.e.* their model is *dense*. Note that our model does not impose dynamic structure on the latent intensities, however, we verified empirically that adding a specification like that of [31] to our model does not improve the performance. We believe that this is the case because our model already allows for capturing complex nonlinear temporal dynamics, thus the addition of linear dynamics to the intensities does not meaningfully impact the model representation abilities.

All prior work on dynamic models using either RBMs, SBNs and Poisson factor analysis assumed uniform temporal sampling. The approach to nonuniform sampling introduced here specifies the dynamics of the transition model in terms of Poisson rates  $\tilde{\boldsymbol{\lambda}}_{nt} = \tau_{nt}^{-1} \Phi(\mathbf{w}_{nt-1} \circ \mathbf{h}_{nt-1}) + \tilde{\boldsymbol{\lambda}}_0$ ; the explicit dependence on sampling delay  $\tau_{nt}$  and scaling  $\mathbf{w}_{nt-1}$  moves well beyond only depending on the  $2^K$  variants of the states  $\mathbf{h}_{nt-1}$ .

#### V. EXPERIMENTS

We present extensive experiments on five datasets. The first consists of artificially generated data with different levels of sparsity, selected to show the computational efficiency of our learning procedure, implemented using GPUs. The next two datasets are publicly available, and consist of collections of text documents over time. Specifically, we consider the US presidential State of the Union addresses [32], and NIPS abstracts [33]. We also consider one public dataset of binary time-series data, consisting of multiple polyphonic music transcriptions [10]. Finally, we consider a microbiome dataset composed of measurements of human gut microbiota over time from 6 subjects spanning 3 different studies [34] (sampled nonuniformly in time). The source code of our dynamic Poisson factor model will be available upon publication.

##### A. Artificial data

We wish to evaluate quantitatively the parallelized implementation of our model using GPUs, to sample from the multinomial distribution in (9), in terms of runtime. We compare this *efficient* implementation with our own Matlab and C++ implementation, which only differs from the GPU version in the multinomial sampler routine, where most of the runtime is spent. There are two routines implemented in C++ in both versions, the CRT and the truncated Poisson samplers. For the experiment we used a standard desktop machine with 4 cores at 3.2GHz and 24Mb RAM. The GPU is a Geforce GTX

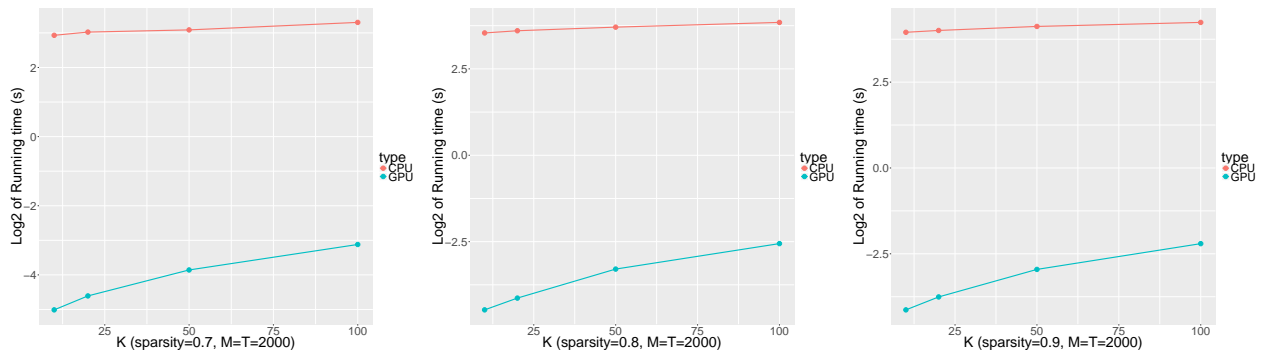


Fig. 2. Computational complexity of dynamic Poisson factor analysis on artificial data.  $\text{Log}_2$ -transformed runtime per full Gibbs iteration (in seconds) vs. number of binary latent variables, for datasets of fixed size, but different sparsity levels.

750i with 640 cores and 2Mb RAM. We consider datasets of size  $M = T = 2000$ ,  $N = 1$  and different observed sparsity levels, namely fractional levels of sparsity equal to members of the set  $\{0.7, 0.8, 0.9\}$ . Figure 2 shows average runtime in seconds for a full Gibbs iteration (a cycle through the updates in Section III), as a function of the number of binary latent variables,  $K$ , in the dynamic Poisson factor model. In all cases we observe speedups ranging from 85 to 250x with an average of about 120x, which constitutes a substantial acceleration, considering the relatively outdated GPU we used for the experiments. It is worth noting that sampling from multinomial distributions represents about 90% of the total runtime.

### B. State of the Union

This dataset contains transcripts of  $T = 225$  US presidential State of the Union addresses, ranging from years 1790 to 2014. We consider the dataset as a single time-series ( $N = 1$ ) with 225 time-points, where each transcript is a document, thus one document per year ( $\tau_{nt} = 1, \forall t$  and  $n = 1$ ). We preprocess the data lightly by removing stop words and terms that occur less than 7 times in one document or less than 20 times overall, which results in a vocabulary of size  $M = 2375$  terms. This preprocessing scheme was previously used by [12] in their experiments with TSBNs.

a) *Quantitative evaluation:* For prediction tasks, we exclude the last year (2014) from the learning phase of the model. For all the other documents ranging from year 1790 to 2013, we randomly partition words from each document into a 80%/20% split. The learning phase of the model is performed on the 80% subset, whereas the remaining 20% observations are used during inference to make predictions at each year. The predictions from both held-out sets are ranked according to their average predicted counts using the emission model in (2). For our model we run 900 iterations of the Gibbs sampler during learning and average predictions over 100 collection samples during inference. We verified empirically that increasing the number of Gibbs iterations does not significantly change results.

To evaluate the prediction performance, we calculate the precision @top- $L$  as in [31], which is defined as the fraction

of the top- $L$  words, predicted by the model, that match the true ranking of the observed word counts. We use  $L = 50$  as in [12]. We compare our dynamic Poisson model against three recently proposed and related models: Gamma Process Dynamic Factor Analysis (GP-DPFA) [31], Dynamic Rank Factor Model (DRFM) [32] and Temporal Sigmoid Belief Network (TSBN) [12]. The parameters of each of these models were selected to maximize performance. For DRFM and TSBN we used 25 latent variables and for GP-DPFA and our model, 100 latent variables. Note that unlike GP-DPFA, our model has latent binary activations that allow for model size selection, thus  $K = 100$  can be seen as a lower bound on the total number of active latent variables. For this dataset we observed that the model estimated non-trivial activations for an average of 24 latent variables; see qualitative evaluation below for further details.

TABLE I  
AVERAGE PREDICTIVE PRECISION FOR STATE OF THE UNION DATASET. MEAN PRECISION WAS COMPUTED ON HELD-OUT SUBSETS OF EACH YEAR IN THE DATASET. PREDICTIVE PRECISION WAS COMPUTED ON THE LAST YEAR, 2014.

Model	Mean Precision	Predictive Precision
Dynamic PFA	0.382	0.520
TSBN	0.327	0.353
GP-DPFA	0.223	0.189
DRFM	0.217	0.177

We summarize predictive performance results in Table I. We report Mean Precision over all documents (years) in the dataset, restricted to the 20% held-out sets, one per year. We also report Predictive Precision for the final year, 2014. We see that our model significantly outperforms all the other methods in terms of mean precision and is comparable to TSBN in terms of predictive precision.

b) *Qualitative evaluation:* We now examine some of the parameters learned by the model to highlight the interpretability of our model. By looking at the conditional posterior of  $r_k$ , the global shape of intensities,  $\theta_{kn}$ , we verified that the model only uses 24 latent variables from the  $K = 100$  set *a priori*, see Figure 3(Top-left). This means that the remaining 76 latent variables are not active,  $h_{knt} = 0$ , or their intensities are close zero,  $\theta_{knt} \approx 0$ , at any given time point. In Figure

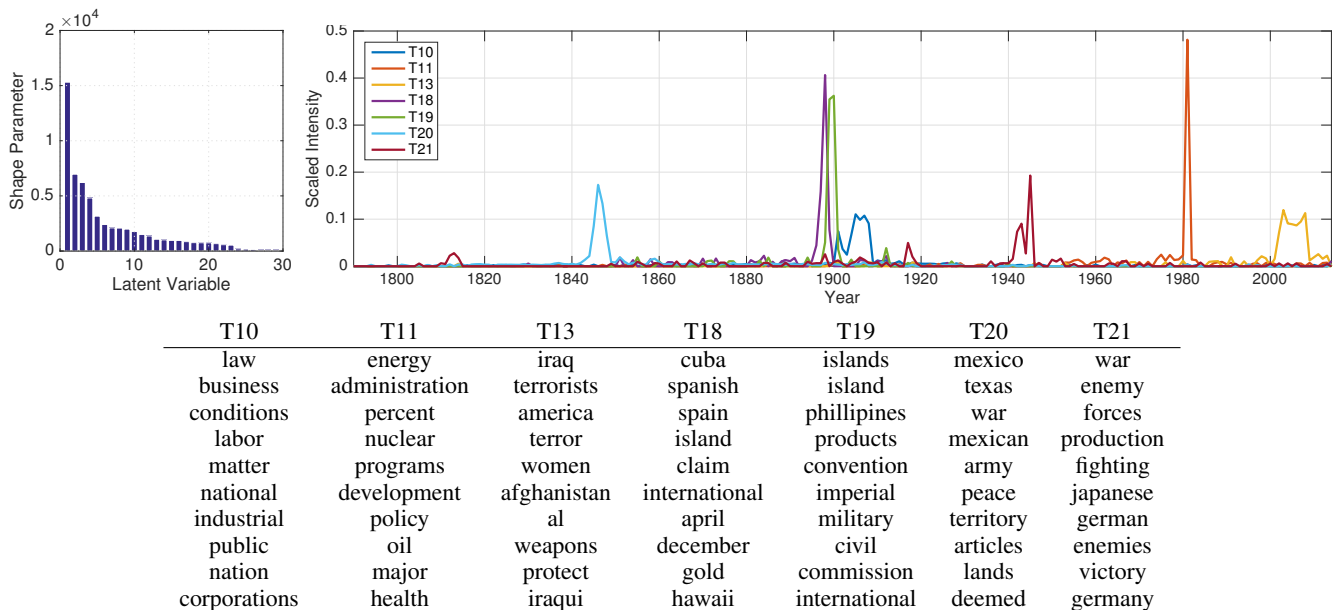


Fig. 3. Selected topics learned from the State of the Union dataset. Top-left: Posterior mean of the shape parameter,  $r_k$ , for latent intensities,  $\theta_{knt}$ . We show the top-30 largest values sorted in decreasing order. Top-right: Scaled intensities,  $\theta_{knt} / \sum_t \theta_{knt}$ , for selected topics. Scaling was done only to improve visualization. Bottom: Top words from selected topics. Each column represents the top-10 largest weights from 7 columns,  $\psi_k$ , of loadings matrix,  $\Psi$ , with corresponding intensity traces shown above.

3(Top-right) and 3(Bottom) we show intensity traces ( $\theta_{knt}, \forall t$ ) and top words (largest weights of  $\psi_k$ ), respectively, from 7 of the 24 active topics estimated by the model. We observe very relevant topics, tightly localized in time. We see for example that Topic 10 is related to the organized labor movement, Topic 11 with the National Energy Program, Topics 13, 18, 19 and 20 focus on the Middle East, Spanish, Phillipines and Mexican wars, respectively. Finally Topic 21 is related to the world wars. Topics not shown are less localized in time and range from foreign policy to internal economic affairs.

### C. NIPS Abstracts

This dataset contains distributions of words (including authors) in all NIPS papers from years 1988 to 2003. Again, we consider the dataset as one time-series ( $N = 1$ ) with 17 time-points, where each year is a document, thus  $\tau_{nt} = 1, \forall t$  and  $n = 1$ . We preprocess the dataset using the same criteria used for the State of the Union data, which results in  $M = 14,036$  distinct words.

TABLE II  
AVERAGE PREDICTIVE PRECISION FOR NIPS ABSTRACTS. MEAN PRECISION WAS COMPUTED ON HELD-OUT SUBSETS OF EACH YEAR. PREDICTIVE PRECISION WAS COMPUTED ON THE LAST YEAR.

Model	Mean Precision	Predictive precision
Dynamic PFA	0.876	0.664
TSBN	0.810	0.538

In this experiment, we focus on a quantitative evaluation using the performance metrics previously defined. Provided that TSBN consistently outperforms DRFM, GP-DPFA and TRBM [12], we only consider comparisons against TSBN

here. In both models, the number of latent variables is set to  $K = 50$ . Predictive performance is summarized in Table II, from which we can see that our model outperforms TSBN by marked margins, using a 80/20% split for mean precisions and the last year, 2003, for predictive precisions.

### D. Music

In this experiment we evaluate the model specification for binary time-series described in Section II-C. For this purpose, we employ the so-called music dataset [10]. This dataset consists of four pieces of polyphonic music sequences of piano, in which time-points, uniformly sampled in time, are encoded as 88-dimensional binary vectors indicating what keys of the piano, ranging from A0 to C8, are “pressed” at any given point in time. The number of simultaneous keys (polyphony) ranges from 0 (silence) to 15 with an average of 3.9 keys. The four pieces in the dataset, namely, Nottingham, Piano, Muse and JSB, correspond to folk tunes, a classical piano MIDI archive, orchestral classical music and harmonized chorales by J.S. Bach, respectively, and span different playing styles and tempo.

Provided that the observed data is binary, in this experiment we report area under the ROC curve (AUC) on the last time-point of each music piece, as performance metric. The last time point was not used during the learning phase. Table III shows once again that our dynamic PFA model consistently outperforms TSBN. It is worth mentioning that [12] reported results on these data, but using marginal log-likelihood estimates as performance metric. We opted for AUCs because we consider this a more fair metric, provided that for TSBN we will have to report a variational lower bound (conservative

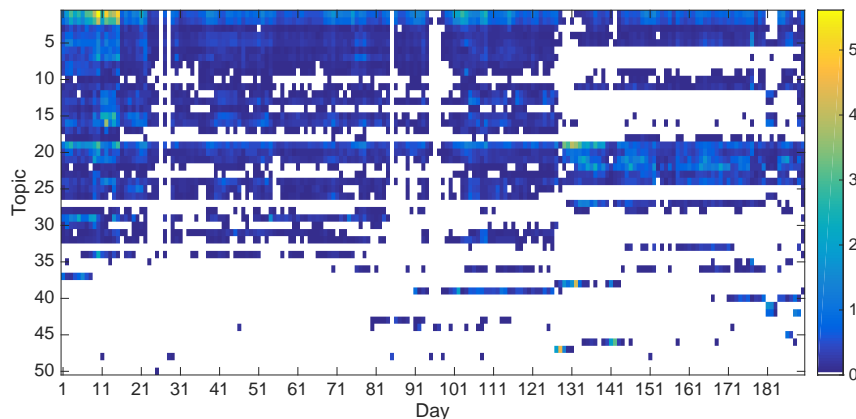


Fig. 4. Intensity heatmap for microbiome data with 50 topics (y axis). The x axis represents time in days. Topics are sorted by usage: top topics (rows) are used throughout the time series whereas bottom topics are highly localized in short time segments. Uncolored topic-day combinations mean that topic  $k$  is not used at a given point time point  $t$ ,  $h_{knt} = 0$ .

estimate) whereas in our case we will have to report a Monte Carlo estimate based on annealed important sampling, which is known to have the potential for overestimating the true marginal log-likelihood.

TABLE III  
AUC FOR MUSIC DATA. PREDICTIVE AUC WAS COMPUTED ON THE LAST TIME POINT OF EACH MUSIC PIECE.

Dataset	Dynamic PFA	TSBN
Piano	0.9303	0.8784
Muse	0.9742	0.9466
JSB	0.9737	0.9686
Nottingham	0.9978	0.9964

### E. Microbiome

We conclude by considering a microbiome dataset composed of longitudinal measurements of human gut microbiota over time, from 6 subjects spanning 3 different studies, with 2 subjects per study. Details about sample collection and processing are detailed in [34]. Data are produced by DNA sequencing of microbiota samples, followed by processing and mapping of raw DNA reads into Operational Taxonomic Units (OTUs). Each OTU defines a species or a group of species, and is commonly used as unit of microbial diversity and is represented in the data as a count, which is a proxy for OTU concentration. The total number of OTUs per subject is shown in Table IV. The sparsity level of these data is on average 85% (most OTUs are not observed at a given point in time), and this sparsity is leveraged in the proposed model (via the BPL link) to yield significant computational acceleration.

Importantly, these data are sampled *nonuniformly* in time, and this complexity motivated several aspects of the proposed model (as detailed when presenting the model). All of the previous models against which we compared in the previous examples are not applicable here, as they assume uniform temporal sampling. Hence, in this experiment we focus on results based upon our proposed model. As in the other experiments, we run 900 Gibbs iterations during learning, 100

collections samples during inference and we set the number of latent variables to  $K = 50$ . Each dataset, has time-series of different time lengths, sampling intervals and OTU resolutions (see Table IV). In particular, data for which  $T = 30$  were sampled once a day, whereas all the others were sampled nonuniformly over a total period of a year [34].

TABLE IV  
ONE-STEP AHEAD FORECASTING RESULTS ON MICROBIOME DATA, IN TERMS OF CORRELATION. ERROR BARS FOR DYNAMIC PFA CORRESPOND TO CORRELATION AVERAGES OVER 100 POSTERIOR COLLECTION SAMPLES DURING INFERENCE.

Sample	$M$	$T$	Dynamic PFA	Naive
S1	5432	321	$0.880 \pm 0.008$	0.8614
S2	5432	189	$0.755 \pm 0.044$	0.378
S3	9371	30	$0.989 \pm 0.003$	0.990
S4	9371	30	$0.964 \pm 0.006$	0.960
S5	33750	332	$0.943 \pm 0.003$	0.935
S6	33750	129	$0.975 \pm 0.002$	0.943

We have highlighted above aspects of the transition-statistics model that are motivated by nonuniform temporal sampling. We here also emphasize components of the emission statistics that are driven by the motivating microbiome example. Recall the emission model,  $\mathbf{x}_{nt} \sim \text{Poisson}(\Psi(\boldsymbol{\theta}_{nt} \circ \mathbf{h}_{nt}))$ , where  $\mathbf{h}_{nt}$  defines which of the  $2^K$  states characterizes time point  $t$  in sample  $n$ , with  $\boldsymbol{\theta}_{nt}$  providing a time- and sample-dependent *scaling* of the factors that are “on” (those having non-zero values in  $\mathbf{h}_{nt}$ ). In microbiome data, the *absolute* counts of OTUs may depend on the sample size and other aspects of the samples, and hence the scaling flexibility provided by  $\boldsymbol{\theta}_{nt}$  is important. This is analogous to analyzing a corpus of documents for which the absolute size of the documents may vary widely.

We also observed that the Bernoulli-Poisson link is particularly useful in our microbiome dataset analysis. One well-known fact in omics studies in general is that less enriched microbiome species may have greater impact to the host, compared with species with high abundance [35]. In many cases, the “*existence*” of a species is more important than



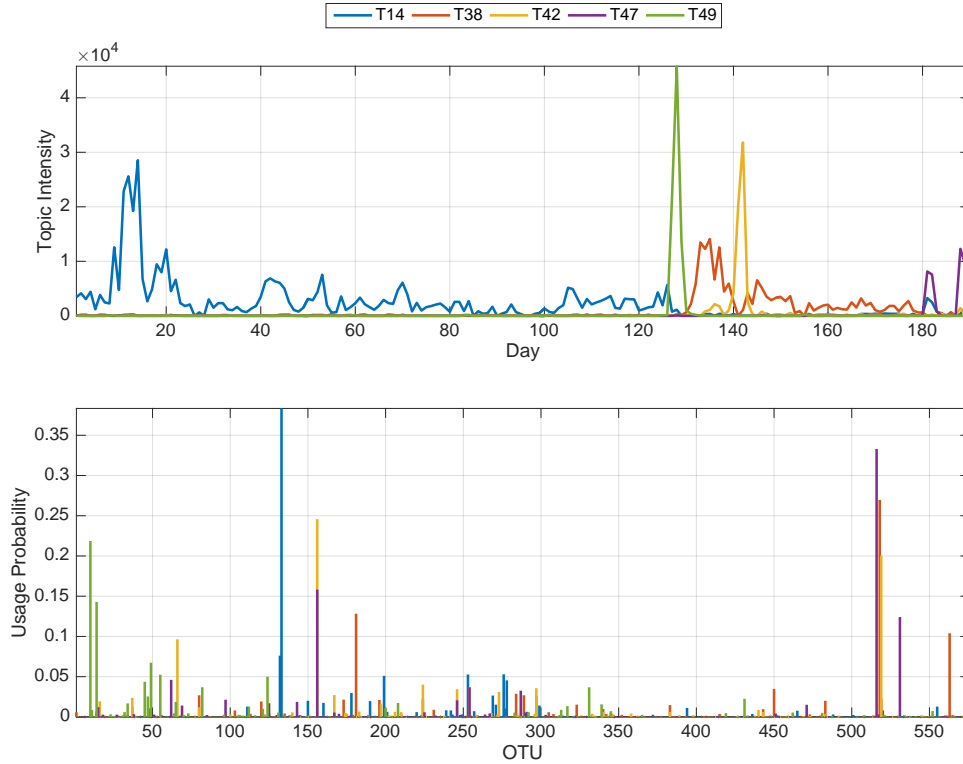


Fig. 5. Selected topics learned from microbiome data and particular to subject S5. Intensities,  $\theta_{nt}$ , and OTU weights,  $\phi_{mk}$ , for selected topics (T14, T38, T42, T47 and T49), in Top and Bottom panels, respectively. The bottom panel of the figure shows weights for 5 different topics (colors), where each bar represent an element of  $\psi_k$ , a column of  $\Psi$ .

the exact value representing its “abundance”. Similar to logit and probit link functions, the non-linearity manifested by BPL latent variables accounts for the saturation effect found in omics data, by diminishing the probability improvement per increment of abundance counts, as counts become large. As a consequence, the model is sensitive to low abundant species while is robust under the scenario where the difference in abundance among species is large. Meanwhile, unlike the logit and probit link, the BPL is asymmetric, which is beneficial when the probability of high-abundant events and low-abundant events has large discrepancies. The probability  $\pi$  increases sharply near  $p = 0$  but slowly near  $p = 1$ , which makes it ideal when the probability of an event is very small. We observed that in our microbiome data, this is usually the case.

We first evaluate whether our model can make predictions at the last time-point, *i.e.*, one-step ahead forecasting, that correlate better with the ground truth measurements, compared to a *naïve* approach in which we make predictions by assuming that observed OTUs at time  $T$  and  $T - 1$  do not change at all, *i.e.*,  $\mathbf{x}_{nT} = \mathbf{x}_{nT-1}$ , which is known to be a good assumption in some cases. Table IV shows Spearman correlations indicating that for 5 out of 6 time-series, modeling the dynamics of OTU changes over time has actual predictive value. We see these results on forecasting of OTU concentrations as interesting

preliminary results that need to be further investigated, thus we leave them as future work.

Figure 4 and 5 show the heatmap of estimated intensities, and five selected topics from the model learned for subject S2, respectively. The proportion of non-zero intensities is 64%. We see that latent variable intensities,  $\theta_{nt}$ , are nicely localized in time. We verified that Topic 49 is consistent with the onset of a *Salmonella* infection suffered by the subject (see [34]), while Topic 38 is related to its recovery period. Interestingly we see that Topic 14, stably present up to the time of infection does not reappear even after recovery has taken place. Also interesting is that the five selected topics in Figure 5 only account for about 10% of the total OTUs in the sample ( $\phi_{mk} > 10^{-4}$ ), which indicates that topics are not only localized in time but in OTU space. In particular, Topic 49 is enriched for *Proteobacteria*, and Topics 14 and 38 are enriched for *Firmicutes*, while Topic 42 is enriched for *Tenericutes*. All these results are consistent with the findings of [34], which were derived using a completely different, more biologically targeted approach.

## VI. CONCLUSION AND FUTURE WORK

We have introduced a dynamic time-series model based on Poisson factor analysis. The model allows for count and binary data, as well as for nonuniformly sampled time-series. Efficient inference using GPUs is developed, that scales with

the number of non-zeros in the data and binary latent variables. Extensive results on benchmark data demonstrate the excellent performance of our simple yet elegant specification. Results on real microbiome data highlight the applicability of our model to interesting problems in modern computational biology.

As future work, we would like to explore the possibility of specifying a deep version of our model, building upon ideas of deep PFA models [15] for extended flexibility and representation ability. We are also considering working on a model specification where multiple observations can occur at the same point in time, but share a common transition *backbone* model. This could be very useful for the analysis of electronic medical records data.

## REFERENCES

- [1] L. Rabiner and B. Juang, "An introduction to hidden Markov models," in *ASSP Magazine, IEEE*, 1986.
- [2] R. Kalman, "Mathematical description of linear dynamical systems," in *J. the Society for Industrial & Applied Mathematics, Series A: Control*, 1963.
- [3] M. Hermans and B. Schrauwen, "Training and analysing deep recurrent neural networks," in *Neural Information Processing Systems*, 2013.
- [4] J. Martens and I. Sutskever, "Learning recurrent neural networks with Hessian-free optimization," in *International Conference of Machine Learning*, 2011.
- [5] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in *International Conference of Machine Learning*, 2013.
- [6] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *International Conference of Machine Learning*, 2013.
- [7] G. Taylor, G. Hinton, and S. Roweis, "Modeling human motion using binary latent variables," in *Neural Information Processing Systems*, 2006.
- [8] I. Sutskever and G. Hinton, "Learning multilevel distributed representations for high-dimensional sequences," in *Artificial Intelligence and Statistics Conference*, 2007.
- [9] I. Sutskever, G. Hinton, and G. Taylor, "The recurrent temporal restricted Boltzmann machine," in *Neural Information Processing Systems*, 2009.
- [10] N. Boulanger-Lewandowski, Y. Bengio, and P. Vincent, "Modeling temporal dependencies in high-dimensional sequences: Application to polyphonic music generation and transcription," in *International Conference of Machine Learning*, 2012.
- [11] R. Mittelman, B. Kuipers, S. Savarese, and H. Lee, "Structured recurrent temporal restricted Boltzmann machines," in *International Conference of Machine Learning*, 2014.
- [12] Z. Gan, C. Li, R. Henao, D. E. Carlson, and L. Carin, "Deep temporal Sigmoid belief networks for sequence modeling," in *Neural Information Processing Systems*, 2015.
- [13] A. Mnih and K. Gregor, "Neural variational inference and learning in belief networks," in *International Conference of Machine Learning*, 2014.
- [14] Z. Gan, C. Chen, R. Henao, D. Carlson, and L. Carin, "Scalable deep Poisson factor analysis for topic modeling," in *International Conference of Machine Learning*, 2015.
- [15] R. Henao, Z. Gan, J. Lu, and L. Carin, "Deep Poisson factor modeling," in *Neural Information Processing Systems*, 2015.
- [16] M. Zhou, "Infinite edge partition models for overlapping community detection and link prediction," in *Artificial Intelligence and Statistics Conference*, 2015.
- [17] M. Zhou, L. Hannah, D. Dunson, and L. Carin, "Beta-negative binomial process and Poisson factor analysis," in *Artificial Intelligence and Statistics Conference*, 2012.
- [18] M. D. Escobar and M. West, "Bayesian density estimation and inference using mixtures," *JASA*, vol. 90, no. 430, pp. 577–588, 1995.
- [19] M. Zhou and L. Carin, "Negative binomial process count and mixture modeling," *Pattern Analysis and Machine Intelligence*, vol. 37, no. 2, pp. 307–320, 2015.
- [20] W. W. Piegorsch, "Complementary log regression for generalized linear models," *The American Statistician*, vol. 46, no. 2, pp. 94–99, 1992.
- [21] D. Collett, *Modelling binary data*. CRC Press, 2002.
- [22] M. Welling and Y. W. Teh, "Bayesian learning via stochastic gradient Langevin dynamics," in *International Conference of Machine Learning*, 2011.
- [23] T. Chen, E. B. Fox, and C. Guestrin, "Stochastic gradient Hamiltonian Monte Carlo," in *International Conference of Machine Learning*, 2014.
- [24] N. Ding, Y. Fang, R. Babbush, C. Chen, R. D. Skeel, and H. Neven, "Bayesian sampling using stochastic gradient thermostats," in *Neural Information Processing Systems*, 2014.
- [25] N. L. Johnson, A. W. Kemp, and S. Kotz, *Univariate discrete distributions*. John Wiley & Sons, 2005, vol. 444.
- [26] G. Taylor and G. Hinton, "Factored conditional restricted Boltzmann machines for modeling motion style," in *International Conference of Machine Learning*, 2009.
- [27] G. Taylor, L. Sigal, D. Fleet, and G. Hinton, "Dynamical binary latent variable models for 3D human pose tracking," in *Conference on Computer Vision and Pattern Recognition*, 2010.
- [28] R. Neal, "Connectionist learning of belief networks," in *Artificial intelligence*, 1992.
- [29] Z. Gan, R. Henao, D. Carlson, and L. Carin, "Learning deep sigmoid belief networks with data augmentation," in *Artificial Intelligence and Statistics Conference*, 2015.
- [30] Z. Gan, C. Chen, R. Henao, D. Carlson, and L. Carin, "Scalable deep poisson factor analysis for topic modeling," in *International Conference of Machine Learning*, 2015.
- [31] A. Acharya, J. Ghosh, and M. Zhou, "Nonparametric Bayesian factor analysis for dynamic count matrices," in *Artificial Intelligence and Statistics Conference*, 2015.
- [32] S. Han, L. Du, E. Salazar, and L. Carin, "Dynamic rank factor model for text streams," in *Neural Information Processing Systems*, 2014.
- [33] A. Globerson, G. Chechik, F. Pereira, and N. Tishby, "Euclidean embedding of co-occurrence data," *Journal of Machine Learning Research*, vol. 8, pp. 2265–2295, 2007.
- [34] L. A. David, A. C. Materna, J. Friedman, M. I. Campos-Baptista, M. C. Blackburn, A. Perrotta, S. E. Erdman, and E. J. Alm, "Host lifestyle affects human microbiota on daily timescales," *Genome Biology*, vol. 15, no. 7, 2014.
- [35] J. Qin, R. Li, J. Raes, M. Arumugam, K. S. Burgdorf, C. Manichanh, T. Nielsen, N. Pons, F. Levenez, T. Yamada *et al.*, "A human gut microbial gene catalogue established by metagenomic sequencing," *nature*, vol. 464, no. 7285, pp. 59–65, 2010.