
Clustering Neural Populations by Poisson Dynamic Factor Model

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Modern recording techniques allow neuroscientists to study multiple neural pop-
2 ulations over extended time periods in a large-scale, and the relationships within
3 and between populations are summarized by low-dimensional latent vectors. When
4 the neural activities are globally nonlinear, using single population analysis is
5 inappropriate. However, defining the populations is usually difficult and wrong
6 cluster assignments will lead to bias in latent structure inferences. To tackle this
7 challenge, we develop a clustering method based mixture of Poisson dynamic
8 factor model. The number of cluster is treated as a parameter in mixture of finite
9 mixtures (MFM) model, and the posteriors are sampled by a MCMC algorithm. To
10 sample the posteriors efficiently, we approximate the full conditional distribution
11 of latent state by Gaussian and approximate the marginal likelihood by making use
12 of the Poisson-Gamma conjugacy. We further apply our method to neuropixel data
13 and hippocampus data for illustration.

14 1 Introduction

15 [TODO]

16 2 Method

17 2.1 Poisson Dynamic Factor Model

18 Denote the observed spike-count of neuron $i \in \{1, \dots, N\}$ at time bin $t \in \{1, \dots, T\}$ as $y_{it} \in \mathbb{Z}_{\geq 0}$,
19 and let $\mathbf{y}_i = (y_{i1}, \dots, y_{iT})'$. Further, let $z_i = j$ denote the cluster indicator of neuron i . To facilitate
20 clustering, we re-parametrize the regular Poisson linear dynamical system (PLDS) model to separate
21 the mean log-firing-rate out. Assume neurons are independently Poisson distributed, conditional on
22 the low-dimensional latent state $\mathbf{x}_t^{(j)} \in \mathbb{R}^p$ as follows:

$$y_{it} \sim \text{Poi}(\lambda_{it})$$
$$\log \lambda_{it} = \mu_t^{(j)} + \mathbf{c}_i' \mathbf{x}_t^{(j)}$$

23 , with $\mathbf{c}_i \sim N(\mathbf{0}, \mathbf{I}_p)$. We further assume the intercept $\mu_t^{(j)}$ and the latent state $\mathbf{x}_t^{(j)}$ progress linearly
 24 with Gaussian noise as

$$\begin{aligned}\mu_1^{(j)} &\sim N(\mu_0, \Sigma_0) \\ \mu_{t+1}^{(j)} &\sim N(f^{(j)}\mu_t^{(j)} + g^{(j)}, \Sigma^{(j)}) \\ \mathbf{x}_1^{(j)} &\sim N(\mathbf{x}_0, \mathbf{Q}_0) \\ \mathbf{x}_{t+1}^{(j)} &\sim N(\mathbf{A}^{(j)}\mathbf{x}_t^{(j)} + \mathbf{b}^{(j)}, \mathbf{Q}^{(j)})\end{aligned}$$

25 If we denote $\boldsymbol{\lambda}_i = (\lambda_{i1}, \dots, \lambda_{iT})'$, $\boldsymbol{\mu}^{(j)} = (\mu_1^{(j)}, \dots, \mu_T^{(j)})'$ and $\mathbf{X}^{(j)} = (\mathbf{x}_1^{(j)}, \dots, \mathbf{x}_T^{(j)})'$, the
 26 model can be equivalently written as regular Poisson factor model as

$$\begin{aligned}\mathbf{y}_i &\sim \text{Poi}(\boldsymbol{\lambda}_i) \\ \log \boldsymbol{\lambda}_i &= \boldsymbol{\mu}^{(j)} + \mathbf{X}^{(j)}\mathbf{c}_i\end{aligned}$$

27 However, since the condition $T/N \rightarrow 0$ doesn't hold [?], the latent state $\mathbf{X}^{(j)}$ cannot be consistently
 28 estimated, and assuming linear dynamics of $\mathbf{X}^{(j)}$ resolves the problem. Note that when $p > 1$, the
 29 model is not unique, since $\tilde{\mathbf{X}}^{(j)} = \mathbf{X}^{(j)}\mathbf{U}$ also satisfies the equation for any orthogonal matrix \mathbf{U}
 30 of order p . To ensure the model indefinability, we simply assume $\mathbf{A}^{(j)}$ and $\mathbf{Q}^{(j)}$ are both diagonal
 31 for convenience. See more detailed discussions of the constraints in discussion. Overall, denote the
 32 cluster-related parameters of cluster j as $\boldsymbol{\theta}^{(j)} = \{\boldsymbol{\mu}^{(j)}, \mathbf{X}^{(j)}, f^{(j)}, g^{(j)}, \Sigma^{(j)}, \{\mathbf{A}^{(j)}, \{\mathbf{b}^{(j)}, \{\mathbf{Q}^{(j)}\}\}$
 33 and the spike counts of neuron i is generated by Poisson dynamic factor model (PDFM) as $\mathbf{Y}_i \sim$
 34 $\text{PDFM}(\boldsymbol{\theta}^{(j)})$, with the prior of $\boldsymbol{\theta}^{(j)}$ as \mathbf{H} . The priors for $\{f^{(j)}, g^{(j)}, \Sigma^{(j)}, \{\mathbf{A}^{(j)}, \{\mathbf{b}^{(j)}, \{\mathbf{Q}^{(j)}\}$
 35 are regular normal and inverse-gamma distribution (or multivariate normal and inverse-Wishart when
 36 using other non-diagonal constraints).

37 The marginal likelihood of neuron i is

$$M_{\boldsymbol{\theta}^{(j)}}(\mathbf{y}_i) = P(\mathbf{y}_i|\boldsymbol{\theta}^{(j)}) = \int P(\mathbf{y}_i|\boldsymbol{\theta}^{(j)}, \mathbf{c}_i)P(\mathbf{c}_i) d\mathbf{c}_i$$

38 The marginal likelihood has no closed form and will be used for clustering. To help with fast clustering,
 39 instead of doing the Laplace approximation, we choose to make use of the Poisson-Gamma conjugacy.
 40 This approximation was originally used in El-Sayyad [1973] to derive approximate posterior and
 41 the same method was applied to derive other approximations in Chan and Vasconcelos [2009].
 42 This approximation leads to the closed form approximation. By the conditional independency
 43 assumption, $M_{\boldsymbol{\theta}^{(j)}}(\mathbf{y}_i) = \prod_{t=1}^T P(y_{it}|\boldsymbol{\theta}^{(j)})$. Since $\mathbf{c}_i \sim N(\mathbf{0}, \mathbf{I}_p)$, $\lambda_{it} = \exp(\mu_t^{(j)} + \mathbf{c}_i'\mathbf{x}_t^{(j)}) \sim$
 44 $\text{lognormal}(\mu_t^{(j)}, \mathbf{x}_t^{(j)'}\mathbf{x}_t^{(j)})$. Approximate the lognormal distribution by Gamma distribution, s.t.
 45 $\text{lognormal}(\mu_t^{(j)}, \mathbf{x}_t^{(j)'}\mathbf{x}_t^{(j)}) \approx \text{Gamma}(a_{it}, b_{it})$ with $a_{it} = (\mathbf{x}_t^{(j)'}\mathbf{x}_t^{(j)})^{-1}$ and $b_{it} = \mathbf{x}_t^{(j)'}\mathbf{x}_t^{(j)} \cdot e^{\mu_t^{(j)}}$.
 46 Then by Poisson-Gamma conjugacy,

$$P(y_{it}|\boldsymbol{\theta}^{(j)}) = \int P(y_{it}|\lambda_{it})P(\lambda_{it}) d\lambda_{it} \approx NB(y_{it}|r_{it}, p_{it})$$

47 , with $r_{it} = a_{it}$ and $p_{it} = 1/(1 + b_{it})$.

48 Another more general idea is to approximate the log-likelihood by second-order polynomials, with
 49 coefficients determined by Chebyshev polynomial approximation Keeley et al. [2019]. However,
 50 the approximation doesn't work well in practice, especially when the neural spike counts have a
 51 wide range. When doing the integration, we need to exponentiate the log-likelihood and this will
 52 exaggerate the approximation error.

53 2.2 Cluster by Mixture of Finite Mixtures Model

54 When the label is unknown, we cluster the neurons by mixture models. In practice, it's usually
 55 impossible to know the number of neural population. One potential method is to do clustering by

56 Dirichlet process mixtures (DPM) model. However, this is conceptually incorrect, since the number
 57 of neural populations is finite but unknown. Besides the conceptual incorrectness, using DPM is not
 58 easy to integrate the field knowledge about the number of neural populations. Here, we choose to use
 59 the mixture of finite mixtures (MFM, Miller and Harrison [2018]) model as follows

$$\begin{aligned}
 K &\sim p_k && \text{where } p_k \text{ is a p.m.f. on } \{1, 2, \dots\} \\
 \boldsymbol{\pi} = (\pi_1, \dots, \pi_k) &\sim \text{Dir}_k(\gamma, \dots, \gamma) && \text{given } K = k \\
 Z_1, \dots, Z_n &\stackrel{i.i.d.}{\sim} \boldsymbol{\pi} && \text{given } \boldsymbol{\pi} \\
 \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k &\stackrel{i.i.d.}{\sim} \boldsymbol{H} && \text{given } K = k \\
 \mathbf{Y}_i = (y_{i1}, \dots, y_{iT})' &\sim \text{SSFM}(\boldsymbol{\theta}^{(z_i)}) && \text{independently for } i = 1, \dots, N, \text{ given } \boldsymbol{\theta}_{1:K} \text{ and } Z_{1:N}
 \end{aligned}$$

60 Besides the conceptual correctness, using MFM model allows us integrate the prior knowledge easily.
 61 Moreover, compared to DPM, MFM has some better properties for clustering, for example, MFM
 62 posterior on number of cluster is more concentrated and consistent, and MFM tend to give clusters
 63 size at the same order of magnitude while DPM may lead to a few large clusters and many small
 64 clusters. See [reference] for detailed discussion.

65 2.3 Inference

66 In this paper, we choose to do inference by MCMC. Because of the Poisson likelihood, the latent
 67 state $\mathbf{X}^{(j)}$ has no closed full conditional distribution. We can sample the posterior by particle MCMC
 68 directly, but this can be slow. However, due to the Markovian structure of the model, the conditional
 69 log-posterior is concave and its Hessian is block-tridiagonal. Thus, we can do the global Laplace
 70 approximation efficiently in $\mathcal{O}(T)$ [Paninski et al., 2010]. The cluster index and number of cluster are
 71 sampled by the analog of partition-based algorithm in DPM [Neal, 2000]. See details of the MCMC
 72 in appendix.

73 In practice, using variational Bayes (VB) instead of MCMC may be more favorable. The PLDS can
 74 be updated by variational EM. Using the stick-breaking representation of MFM model, we can do VB
 75 easily similar to Blei and Jordan [2006]. However, checking by the “gold standard” MCMC before
 76 doing VB is always a good choice.

77 For the dimensionality p of the latent state, we can treat it as a parameter and sample the posterior by
 78 RJMCMC as in Lopes and West [2004] or borrow the idea of adaptive Gibbs sampling with shrinkage
 79 prior [Bhattacharya and Dunson, 2011]. Here, we simply pre-set the p or select the optimized value
 80 by the cross-validation, which can be easily conducted when switching to the deterministic algorithm
 81 in the future.

82 3 Simulations

83 3.1 Model Global Non-linearity by Clustering

84 There were a rich research results for single PLDS model, but it provides only a global linear model
 85 to represent the data in a lower dimensional subspace, which makes the application scope limited.
 86 When the input space is nonhomogeneous, a large dimension of latent state is needed and this may
 87 lead to the overfitting and poor performance. Mixture of PLDS/ PDFM models allow us to partition
 88 the input space into clusters and can therefore capture global nonlinearity by combining local linear
 89 models.

90 To show the connection between the proposed model and PLDS parametrization, we simulate the data
 91 using the PLDS model, although these two are equivalent after some algebra. Denote the common
 92 PLDS model as $\log \lambda_i = \delta_i \mathbf{1}_T + \mathbf{G}^{(j)} \mathbf{d}_i$, where $\mathbf{d}_i \in \mathbb{R}^q$. In PLDS parametrization, the latent state
 93 has one more dimension, i.e. $p = q + 1$.

We first simulated three neural populations, with 10 neurons in each. We set $q = 2$ for each cluster and the recording length is $T=1000$. After checking the trace plots up to 10,000 iterations, the convergence achieved after several steps. The panel A and B in figure 1 show the posterior mean firing rate and fitted latent state, averaging from iteration 500 to 1000. The results are transformed to PLDS parametrization for comparison. Notice that the PLDS parametrization also doesn't have the unique solution, we let $\mathbf{G}^{(j)}$ has mean zeros columns and is orthogonal.

Still has the sign flipping issue. Fix later.

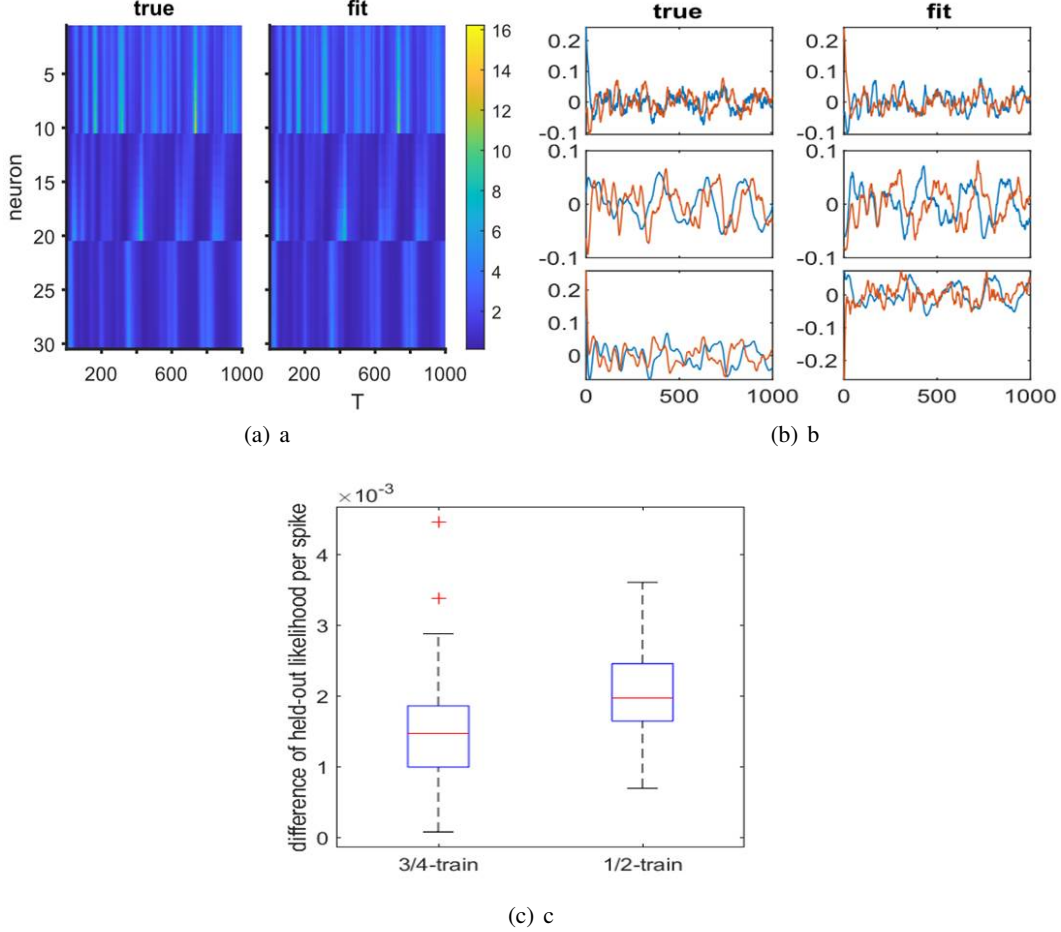


Figure 1: Model Global Non-linearity by Clustering

Then we held out 1/4 and 1/2 data as test set in a “speckled” pattern, i.e. randomly select subset of data for each neuron as held-out dataset. Then we fit the model with and without clusters, keeping the same latent dimension, i.e. (1) 3 clusters with $p = 1$ for each and (2) 1 cluster with $p = 2 \cdot 3 - 1 = 5$. The procedure is replicated for 100 times for two proportions, and the difference of held-out likelihood per spike between 2 fittings are shown in panel C in figure 1. The difference between (1) and (2) are always positive, and this shows doing the mixture of PDFM performs better than single PDFM. Moreover, as the proportion of training decrease (less data), the benefit of clustering becomes more significant. This suggests that doing clustering is necessary.

3.1.1 Clustering

Then use the same setting, we remove the label and use the mixture model to do the clustering by MFM. The prior of the cluster number is $K \sim \text{Geometric}(0.2)$. The panel A of figure 2 shows the trace of label z_i for each neuron for the first 100 iterations.

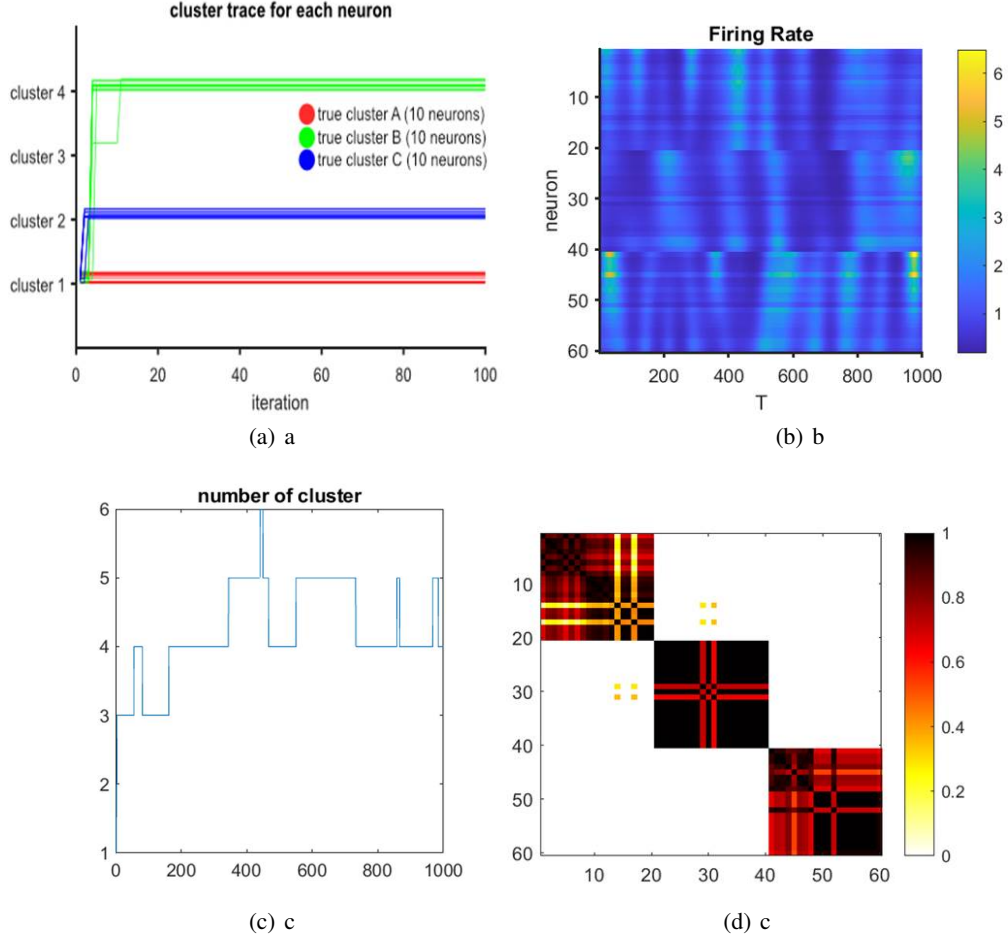


Figure 2: Clustering

In this toy simulation example, the signal is strong and the spiking amplitudes of neurons in each population are similar. This might be too simplified for real situation. In practice, the neural activity is usually sparse and can be recorded in high resolution. Moreover, the spiking amplitude for neuron in one population vary a lot. Therefore, we simulate another more realistic example. In this simulation, there are 3 clusters and the dimension of latent state is $p = 1$ for each. However, there are 20 neurons with wider range but smaller value of loading c_i . The simulated firing rate is shown in panel B in figure 2. The trace plot of number of cluster (panel C in figure 2) shows that the model tend to further split some clusters into sub-population, based on the spiking amplitude. The similarity matrix (panel D in figure 2) of posterior, averaging from iteration 500 to 1000, shows that the model tend to split cluster 1 and 3 into two sub-clusters.

4 Application

4.1 Neuropixels data

[TODO]: brief introduction of Neuropixels dataset.

The Neuropixels dataset contains recording of neural activities in different brain regions, when doing the drift-grating experiment. Here, we only consider neurons with $\text{SNR} > 3$ and the spiking counts > 1000 . Then, we use the recording activity from Lateral posterior nucleus of the thalamus (LP, 20 neurons), anteromedial visual area (VISam, 12 neurons) and ventral posteromedial nucleus of the

thalamus (VPM, 14 neurons) during the spontaneous period. The bin size is 0.1s and truncate the data up to 1000 steps (100s recording) for convenience.

Panel A in figure 3 shows the spiking counts of these 46 neurons. Here we first fit model using all data, with $p = 1$ and $K \sim \text{Geometric}(0.3)$. For the number of clusters, panel B in figure 2 shows the trace plot for the first 5000 iterations and panel C shows the histogram of iteration 2500 to 5000. These plots show that these neurons are quite non-homogenous and they tend to form many sub-populations. The posterior similarity matrix, averaging from iteration 2500 to 5000 (panel D in figure 3), shows that there are several sub-populations in each anatomical cluster, and there are some tiny confusions between VISam and VPM.

I tried to held-out half of the data in a speckled pattern and fit the model with clustering on and off (single population). The dimension is selected by held-out likelihood, which is $p=1$ for clustering on and $p=2$ for single population analysis. The trace plots of the held-out log-likelihood per spike (starting from iteration 2) shows that doing clustering doesn't improve things... Maybe because "turning clustering on" introduce too much variance... The benefit of clustering will pop out if we use the deterministic algorithm, such as VB mentioned.

4.2 Hippocampus data

What kind of story do I need to tell here? This should be different from the story from Neuropixel data. Maybe focus on the potential improvement of held-out likelihood after clustering? But I guess similar problem will happen. When fitting data with 84 neurons, the number of component jumps around 10.

5 Discussion

As mentioned above, the factor model doesn't have unique solution. Although this issue had been thoroughly discussed in statistics [reference], it has been ignored in some neuroscience research. Since the (P)LDS related model are usually fitted by deterministic algorithms such as variants of EM, the issue is kind of hard to detect. And because of the ignorance, there are some problematic comments on linear dynamics \mathbf{A} and \mathbf{Q} . When using the variants of factor model, e.g. (P)LDS and GPFA, we should only focus on the "shape" of latent state but not the specific values of them.

In this paper, to ensure the unique solution, we put the diagonal constraints in $\mathbf{A}^{(j)}$ and $\mathbf{Q}^{(j)}$ are used for convenience. However, only constraining on $p(p-1)/2$ is enough [reference]. Since we the label is unknown and will be switched when doing clustering, it's inappropriate and convenient to put constraint on c_i . Therefore, we instead treat $\mathbf{X}^{(j)}$ as the "loading" and can put constraints on it. In traditional Gaussian factor model, there are two equivalent types of constraints: (1) diagonal constraint: $\mathbf{X}'^{(j)}\mathbf{X}^{(j)}$ is diagonal; (2) block lower triangular constraint [Fokoué and Titterton, 2003]: The first p rows of $\mathbf{X}^{(j)}$ is lower triangular and the diagonal elements are positive. Since we assume $\mathbf{X}^{(j)}$ has linear progression, using the diagonal $\mathbf{X}'^{(j)}\mathbf{X}^{(j)}$ constraint is more appropriate. This constraint can be easily achieved when conducting deterministic optimization algorithms (e.g. the VB mentioned in "method-inference" section). This can also be done in MCMC, by turning off the reflection of the projection (e.g. let the diagonal elements in the projection matrix be positive). Using diagonal $\mathbf{X}'^{(j)}\mathbf{X}^{(j)}$ leads to similar results.

Moreover, the idea of doing clustering by mixture model can be extended beyond the Poisson distribution. We can further include other information, such as, the dispersion information by assuming negative-binomial distributed or even Conway-Maxwell-Poisson distributed neural spikes [reference to our paper in CMP]. This further information can be used for more detailed clustering, by expanding the state space.

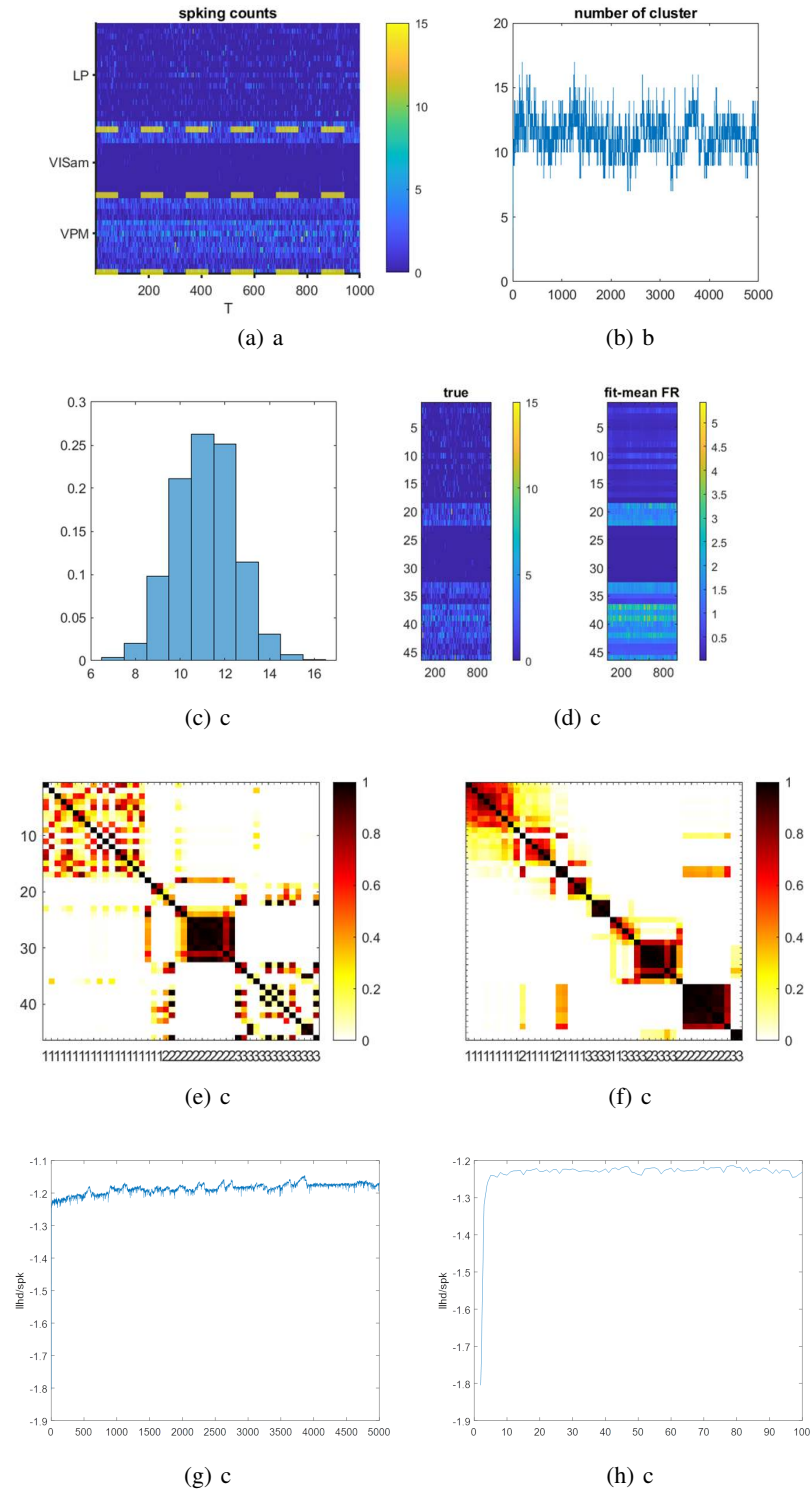


Figure 3: Neuropixel

174 **Acknowledgments**

175 **Broader Impact**

References

- G. M. El-Sayyad. Bayesian and Classical Analysis of Poisson Regression. *Journal of the Royal Statistical Society: Series B (Methodological)*, 35(3):445–451, jul 1973. ISSN 2517-6161. doi: 10.1111/J.2517-6161.1973.TB00972.X.
- Antoni B. Chan and Nuno Vasconcelos. Bayesian poisson regression for crowd counting. *Proceedings of the IEEE International Conference on Computer Vision*, pages 545–551, 2009. doi: 10.1109/ICCV.2009.5459191.
- Stephen L. Keeley, David M. Zoltowski, Yiyi Yu, Jacob L. Yates, Spencer L. Smith, and Jonathan W. Pillow. Efficient non-conjugate Gaussian process factor models for spike count data using polynomial approximations. *37th International Conference on Machine Learning, ICML 2020*, PartF16814:5133–5142, jun 2019.
- Jeffrey W. Miller and Matthew T. Harrison. Mixture models with a prior on the number of components. *Journal of the American Statistical Association*, 113(521):340, jan 2018. doi: 10.1080/01621459.2016.1255636.
- Liam Paninski, Yashar Ahmadian, Daniel Gil Ferreira, Shinsuke Koyama, Kamiar Rahnama Rad, Michael Vidne, Joshua Vogelstein, and Wei Wu. A new look at state-space models for neural data, aug 2010. ISSN 09295313. URL <https://link.springer.com/article/10.1007/s10827-009-0179-x>.
- Radford M Neal. Markov Chain Sampling Methods for Dirichlet Process Mixture Models. *Journal of Computational and Graphical Statistics*, 9(2):249–265, 2000.
- David M Blei and Michael I Jordan. Variational Inference for Dirichlet Process Mixtures. *Bayesian Analysis*, 1(1):121–144, 2006. URL <http://www.cs.berkeley.edu/~simblei/>.
- Hedibert Freitas Lopes and Mike West. Bayesian Model Assessment in Factor Analysis. *Statistica Sinica*, 14: 41–67, 2004.
- A. Bhattacharya and D. B. Dunson. Sparse Bayesian infinite factor models. *Biometrika*, 98(2):291, jun 2011. ISSN 00063444. doi: 10.1093/BIOMET/ASR013.
- Ernest Fokoué and D. M. Titterton. Mixtures of factor analysers. Bayesian estimation and inference by stochastic simulation. *Machine Learning*, 50(1-2):73–94, jan 2003. ISSN 08856125. doi: 10.1023/A:1020297828025.

A Appendix