
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 [Macke
21 et al., 2011] to separate the mean log-firing-rate out. Assume neurons are independently Poisson
22 distributed, conditional on 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

$$\begin{aligned}\mu_1^{(j)} &\sim N(\mu_0, \Sigma_0) \\ \mu_{t+1}^{(j)} &\sim N(f^{(j)}\mu_t^{(j)} + g^{(j)}, \sigma^{2(j)}) \\ \mathbf{x}_1^{(j)} &\sim N(\mathbf{x}_0, \mathbf{Q}_0) \\ \mathbf{x}_{t+1}^{(j)} &\sim N(\mathbf{A}^{(j)}\mathbf{x}_t^{(j)} + \mathbf{b}^{(j)}, \mathbf{Q}^{(j)})\end{aligned}$$

25 The bias term ($\mathbf{b}^{(j)}$) in the latent state is constant across time, unlike ind PLDS, which is $\mathbf{b}_t^{(j)}$.
 26 Allowing the time-varying bias term will lead to infinite solutions of linear dynamics. If the $\mathbf{A}^{(j)}$ is
 27 one solution, then $\mathbf{A}^{*(j)} = \mathbf{A}^{(j)} + \mathbf{M}$ and $\mathbf{b}_t^{*(j)} = \mathbf{b}_t^{(j)} - \mathbf{M}\mathbf{x}_t^{(j)}$ will be another set of solution, for
 28 any $\mathbf{M} \in \mathbb{R}^{p \times p}$. The same reason for using time constant $g^{(j)}$.

29 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
 30 model can be equivalently written as regular Poisson factor model as

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

31 However, since the condition $T/N \rightarrow 0$ doesn't hold [Johnstone and Lu, 2009], the latent state
 32 $\mathbf{X}^{(j)}$ cannot be consistently estimated, and assuming linear dynamics of $\mathbf{X}^{(j)}$ resolves the prob-
 33 lem. Note that when $p > 1$, the model is not unique, since $\mathbf{X}^{*(j)} = \mathbf{X}^{(j)}\mathbf{U}$ also satisfies the
 34 equation for any orthogonal matrix \mathbf{U} of order p . To ensure the model indefinability, we sim-
 35 ply assume $\mathbf{A}^{(j)}$ and $\mathbf{Q}^{(j)}$ are both diagonal for convenience. See more detailed discussions
 36 of the constraints in discussion. Overall, denote the cluster-related parameters of cluster j as
 37 $\boldsymbol{\theta}^{(j)} = \{\boldsymbol{\mu}^{(j)}, \mathbf{X}^{(j)}, f^{(j)}, g^{(j)}, \sigma^{2(j)}, \mathbf{A}^{(j)}, \mathbf{b}^{(j)}, \mathbf{Q}^{(j)}\}$ and the spike counts of neuron i is generated
 38 by Poisson dynamic factor model (PDFM) as $\mathbf{Y}_i \sim PDFM(\boldsymbol{\theta}^{(j)})$, with the prior of $\boldsymbol{\theta}^{(j)}$ as \mathbf{H} . The
 39 priors for $\{f^{(j)}, g^{(j)}, \sigma^{2(j)}, \mathbf{A}^{(j)}, \mathbf{b}^{(j)}, \mathbf{Q}^{(j)}\}$ are regular normal and inverse-gamma distribution (or
 40 multivariate normal and inverse-Wishart when using other non-diagonal constraints).

41 The marginal likelihood of neuron i is

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

42 The marginal likelihood has no closed form and will be used for clustering. To help with fast clustering,
 43 instead of doing the Laplace approximation, we choose to make use of the Poisson-Gamma conjugacy.
 44 This approximation was originally used in El-Sayyad [1973] to derive approximate posterior and
 45 the same method was applied to derive other approximations in Chan and Vasconcelos [2009].
 46 This approximation leads to the closed form approximation. By the conditional independency
 47 assumption, $M_{\boldsymbol{\theta}^{(j)}}(\mathbf{y}_i) = \prod_{t=1}^T P(y_{it}|\boldsymbol{\theta}^{(j)})$. Since $\mathbf{c}_i \sim N(\mathbf{0}, \mathbf{I}_p)$, $\lambda_{it} = \exp(\mu_t^{(j)} + \mathbf{c}_i'\mathbf{x}_t^{(j)}) \sim$
 48 $lognormal(\mu_t^{(j)}, \mathbf{x}_t^{(j)'}\mathbf{x}_t^{(j)})$. Approximate the lognormal distribution by Gamma distribution, s.t.
 49 $lognormal(\mu_t^{(j)}, \mathbf{x}_t^{(j)'}\mathbf{x}_t^{(j)}) \approx Gamma(a_{it}, b_{it})$ with $a_{it} = (\mathbf{x}_t^{(j)'}\mathbf{x}_t^{(j)})^{-1}$ and $b_{it} = \mathbf{x}_t^{(j)'}\mathbf{x}_t^{(j)} \cdot e^{\mu_t^{(j)}}$.
 50 Then by Poisson-Gamma conjugacy,

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

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

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

57 2.2 Cluster by Mixture of Finite Mixtures Model

58 When the label is unknown, we cluster the neurons by mixture models. In practice, it's usually
 59 impossible to know the number of neural population. One potential method is to do clustering by
 60 Dirichlet process mixtures (DPM) model. However, this is conceptually incorrect, since the number
 61 of neural populations is finite but unknown. Besides the conceptual incorrectness, using DPM is not
 62 easy to integrate the field knowledge about the number of neural populations. Here, we choose to use
 63 the mixture of finite mixtures (MFM, Miller and Harrison [2018]) model as follows

$$\begin{aligned}
 K &\sim p_k && \text{where } p_k \text{ is a p.m.f. on } \{1, 2, \dots\} \\
 \pi = (\pi_1, \dots, \pi_k) &\sim \text{Dir}_k(\gamma, \dots, \gamma) && \text{given } K = k \\
 Z_1, \dots, Z_n &\stackrel{i.i.d.}{\sim} \pi && \text{given } \pi \\
 \theta_1, \dots, \theta_k &\stackrel{i.i.d.}{\sim} \mathbf{H} && \text{given } K = k \\
 \mathbf{Y}_i = (y_{i1}, \dots, y_{iT})' &\sim \text{SSFM}(\theta^{(z_i)}) && \text{independently for } i = 1, \dots, N, \text{ given } \theta_{1:K} \text{ and } Z_{1:N}
 \end{aligned}$$

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

69 2.3 Inference

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

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

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

86 3 Simulations

87 3.1 Model Global Non-linearity by Clustering

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

94 To show the connection between the proposed model and PLDS parametrization, we simulate the data
 95 using the PLDS model, although these two are equivalent after some algebra. Denote the common
 96 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
 97 has one more dimension, i.e. $p = q + 1$.

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

104 **Still has the sign flipping issue. Fix later.**

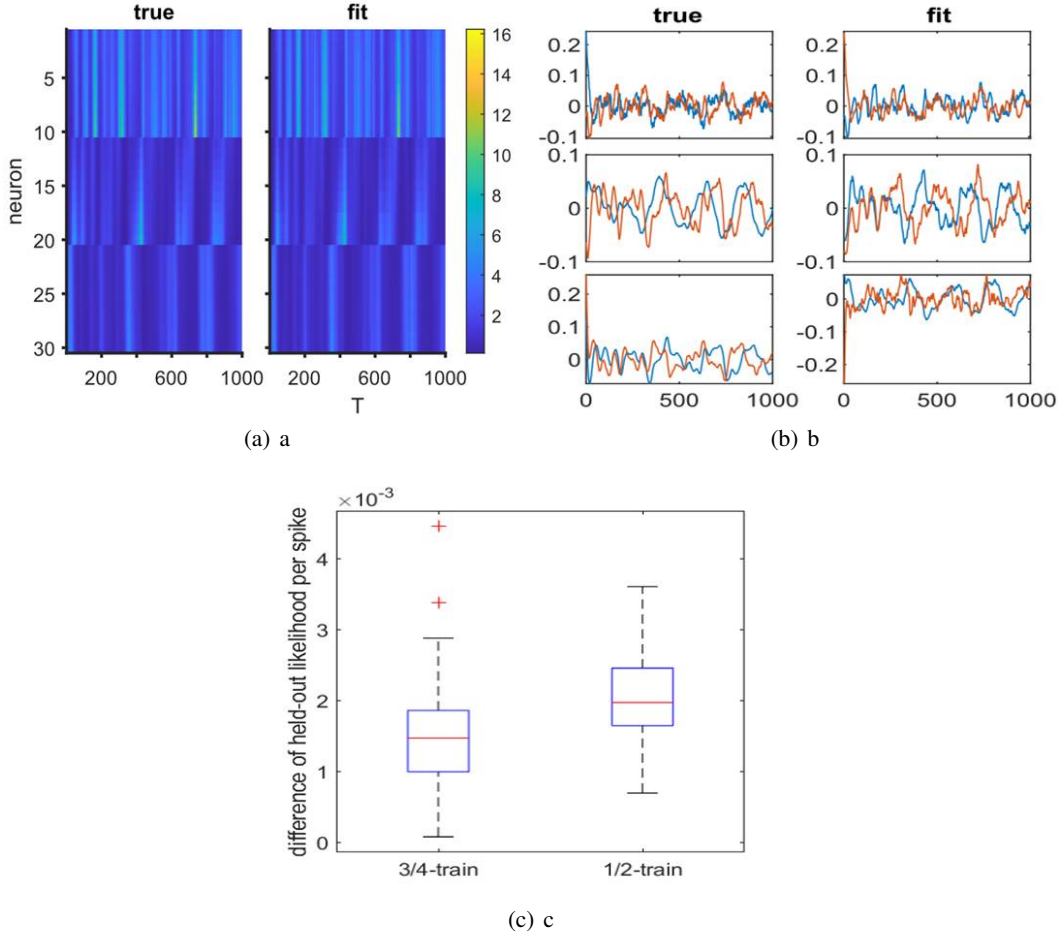


Figure 1: Model Global Non-linearity by Clustering

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

113 3.1.1 Clustering

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

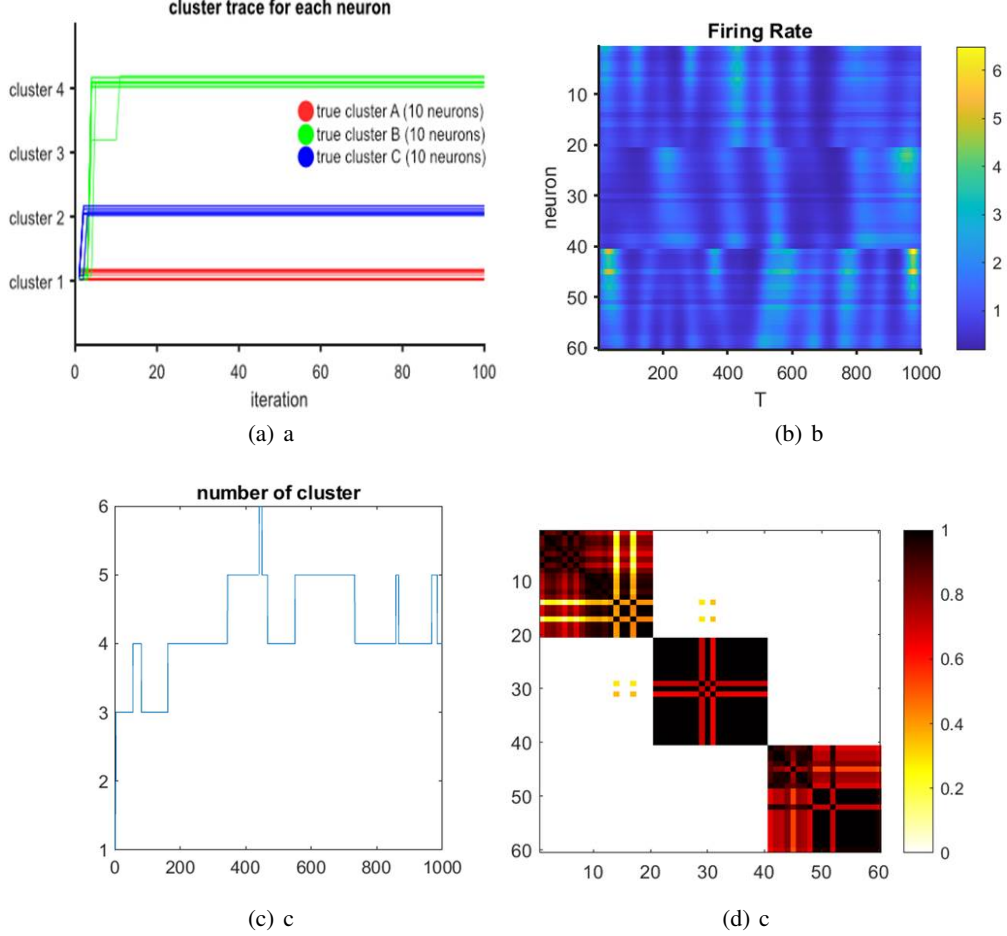


Figure 2: Clustering

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

127 4 Application

128 4.1 Neuropixels data

129 **[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: $\mathbf{X}^{(j)\top} \mathbf{X}^{(j)}$ 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 $\mathbf{X}^{(j)\top} \mathbf{X}^{(j)}$ 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 $\mathbf{X}^{(j)\top} \mathbf{X}^{(j)}$ 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.

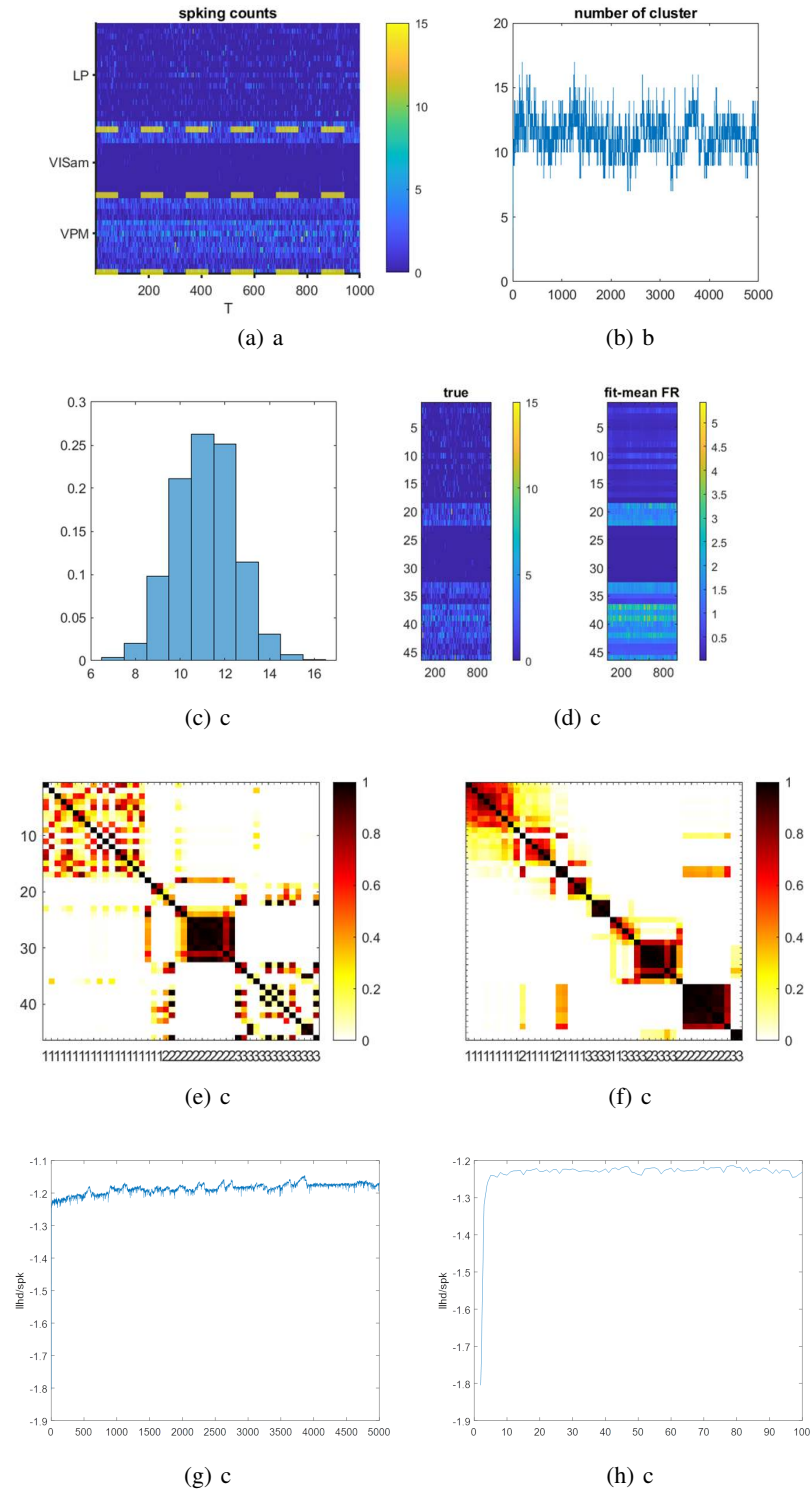


Figure 3: Neuropixel

178 **Acknowledgments**

179 **Broader Impact**

References

- Jakob H. Macke, Lars Buesing, John P. Cunningham, Byron M. Yu, Krishna V. Shenoy, and Maneesh Sahani. Empirical models of spiking in neural populations. *Advances in Neural Information Processing Systems*, 24, 2011.
- Iain M. Johnstone and Arthur Yu Lu. On Consistency and Sparsity for Principal Components Analysis in High Dimensions. *Journal of the American Statistical Association*, 104(486):682, jun 2009. ISSN 01621459. doi: 10.1198/JASA.2009.0121.
- 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/~dblei/>.
- 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.
- Uri T Eden, Loren M. Frank, Riccardo Barbieri, Victor Solo, and Emery N. Brown. Dynamic Analysis of Neural Encoding by Point Process Adaptive Filtering. *Neural Computation*, 16(5):971–998, may 2004. ISSN 0899-7667. doi: 10.1162/089976604773135069. URL <https://doi.org/10.1162/089976604773135069>.
- H E Rauch, F Tung, and C T Striebel. Maximum likelihood estimates of linear dynamic systems. *AIAA Journal*, 3(8):1445–1450, aug 1965. ISSN 0001-1452. doi: 10.2514/3.3166. URL <https://doi.org/10.2514/3.3166>.
- Matthew D. Hoffman and Andrew Gelman. The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research*, 15:1593–1623, nov 2011. ISSN 15337928. URL <https://arxiv.org/abs/1111.4246v1>.

A Appendix

A.1 Full Conditional Distributions

$\mathbf{x}_t^{(j)}$ and $\mu_t^{(j)}$: Let the t^{th} column observation and firing rate of the j^{th} cluster be $\tilde{\mathbf{y}}_t = \text{vec}(\{y_{it}^{(j)} | z_i = j\})$ and $\tilde{\boldsymbol{\lambda}}_t^{(j)} = \text{vec}(\{\lambda_{it} | z_i = j\})$. The number of neurons in cluster j is $n_j = |\{i : z_i = j\}|$. The corresponding loading for these n_j neurons is $\mathbf{C}^{(j)} \in \mathbb{R}^{n_j \times p}$, such that $\log \tilde{\boldsymbol{\lambda}}_t^{(j)} = \mu_t^{(j)} \mathbf{1}_{n_j} + \mathbf{C}^{(j)} \mathbf{x}_t^{(j)} = (\mathbf{1}_{n_j}, \mathbf{C}^{(j)})' \cdot (\mu_t^{(j)}, \mathbf{x}_t^{(j)})'$. Denote $\tilde{\mathbf{C}}^{(j)} =$

226 $(\mathbf{1}_{n_j}, \mathbf{C}^{(j)}), \tilde{\mathbf{x}}_t^{(j)} = (\mu_t^{(j)}, \mathbf{x}_t'^{(j)})', \tilde{\mathbf{x}}^{(j)} = (\tilde{\mathbf{x}}_1^{(j)}, \dots, \tilde{\mathbf{x}}_T^{(j)})', \tilde{\mathbf{A}}^{(j)} = \text{diag}(f^{(j)}, \mathbf{A}^{(j)}), \tilde{\mathbf{b}}^{(j)} =$
 227 $(g^{(j)}, \mathbf{b}'^{(j)})'$ and $\tilde{\mathbf{Q}}^{(j)} = \text{diag}(\sigma^{2(j)}, \mathbf{Q}^{(j)})$. The full conditional distribution $P(\tilde{\mathbf{x}}^{(j)} | \dots) =$
 228 $P(\tilde{\mathbf{x}}^{(j)} | \{\tilde{\mathbf{y}}_t\}_{t=1}^T, \tilde{\mathbf{C}}^{(j)}, \tilde{\mathbf{A}}^{(j)}, \tilde{\mathbf{b}}^{(j)}, \tilde{\mathbf{Q}}^{(j)})$ is approximated by a global Laplace approximation,
 229 i.e. $P(\tilde{\mathbf{x}}^{(j)} | \dots) \approx N_{(p+1)T}(\tilde{\mathbf{x}}^{(j)} | \boldsymbol{\mu}_{\tilde{\mathbf{x}}^{(j)}}, \boldsymbol{\Sigma}_{\tilde{\mathbf{x}}^{(j)}})$, with $\boldsymbol{\mu}_{\tilde{\mathbf{x}}^{(j)}} = \arg \max_{\tilde{\mathbf{x}}^{(j)}} P(\tilde{\mathbf{x}}^{(j)} | \dots)$ and
 230 $\boldsymbol{\Sigma}_{\tilde{\mathbf{x}}^{(j)}} = -(\nabla \nabla \log P(\tilde{\mathbf{x}}^{(j)} | \dots) |_{\tilde{\mathbf{x}}^{(j)} = \boldsymbol{\mu}_{\tilde{\mathbf{x}}^{(j)}}})^{-1}$. The log full conditional distribution $h(\tilde{\mathbf{x}}^{(j)}) =$
 231 $\log P(\tilde{\mathbf{x}}^{(j)} | \dots)$ is given by:

$$h = \text{const} + \sum_{t=1}^T \tilde{\mathbf{y}}_t' \tilde{\mathbf{C}}^{(j)} \tilde{\mathbf{x}}_t^{(j)} - \tilde{\lambda}_t - \frac{1}{2} (\tilde{\mathbf{x}}_1^{(j)} - \tilde{\mathbf{x}}_0^{(j)})' \tilde{\mathbf{Q}}_0^{(j)} (\tilde{\mathbf{x}}_1^{(j)} - \tilde{\mathbf{x}}_0^{(j)})$$

$$- \sum_{t=2}^T \frac{1}{2} (\tilde{\mathbf{x}}_t^{(j)} - \tilde{\mathbf{A}}^{(j)} \tilde{\mathbf{x}}_{t-1}^{(j)} - \tilde{\mathbf{b}}^{(j)})' \tilde{\mathbf{Q}}^{(j)} (\tilde{\mathbf{x}}_t^{(j)} - \tilde{\mathbf{A}}^{(j)} \tilde{\mathbf{x}}_{t-1}^{(j)} - \tilde{\mathbf{b}}^{(j)})$$

232 The $h(\tilde{\mathbf{x}}^{(j)})$ is concave and hence unimodal, and the Markovian structure of the latent dynamics,
 233 and hence the tri-block-diagonal Hessian, makes it possible to compute a Newton update in $\mathcal{O}(T)$
 234 [Paninski et al., 2010].

235 To facilitate convergence, we use a smoothing estimate with local Gaussian approximation as a "warm
 236 start". The forward filtering for a dynamic Poisson model has been previously described in Eden et al.
 237 [2004]. Because the approximated filtering estimates are Gaussian distributed, we can further find the
 238 smoothing estimates using a backward pass as in Rauch et al. [1965].

239 $\underline{\mathbf{c}}_i$: When writing the observation mapping in matrix form, $\log \lambda_i = \boldsymbol{\mu}^{(j)} + \mathbf{X}^{(j)} \mathbf{c}_i$, the update of \mathbf{c}_i
 240 reduces to a regular Poisson regression problem, given $\boldsymbol{\mu}^{(j)}$ and $\mathbf{X}^{(j)}$ are known. Here, we sample
 241 the posterior by no a No U-Turn Sampler (NUTS, Hoffman and Gelman [2011]).

242 Linear dynamics of latent state: The parameters for linear dynamics are $f^{(j)}, g^{(j)}, \sigma^{2(j)}, \mathbf{A}^{(j)},$
 243 $\mathbf{b}^{(j)}$ and $\mathbf{Q}^{(j)}$. To make the model identifiable, we simply assume $\mathbf{A}^{(j)} = \text{diag}(a_1^{(j)}, \dots, a_p^{(j)})$
 244 and $\mathbf{Q}^{(j)} = \text{diag}(q_1^{(j)}, \dots, q_p^{(j)})$. Therefore, we can update $\mathbf{A}^{(j)}, \mathbf{b}^{(j)}$ and $\mathbf{Q}^{(j)}$ dimension-by-
 245 dimension, as the update in $f^{(j)}, g^{(j)}$ and $\sigma^{2(j)}$. Use the priors $\sigma^{2(j)} \sim IG(\nu_0/2, \nu_0 \sigma_0^2/2)$
 246 and $(g^{(j)}, f^{(j)})' \sim N(\boldsymbol{\mu}_0, \sigma^{2(j)} \boldsymbol{\Lambda}_0^{-1})$. The prior $\boldsymbol{\mu}_0$ is specified as $(0, 1)'$. Denote $\mathbf{y}_{\mu^{(j)}} =$
 247 $(\mu_2^{(j)}, \dots, \mu_T^{(j)})'$ and $\mathbf{X}_{\mu^{(j)}} = (\mathbf{1}_{T-1}, \boldsymbol{\mu}_{1:(T-1)}^{(j)})$, with $\boldsymbol{\mu}_{1:(T-1)}^{(j)} = (\mu_1^{(j)}, \dots, \mu_{T-1}^{(j)})'$. The
 248 full conditional distributions are $\sigma^{2(j)} | \dots \sim IG\left(\frac{\nu_0 + T - 1}{2}, \frac{\nu_0 \sigma_0^2 + \mathbf{y}_{\mu^{(j)}}' \mathbf{y}_{\mu^{(j)}} + \boldsymbol{\mu}_0' \boldsymbol{\Lambda}_0 \boldsymbol{\mu}_0 - \boldsymbol{\mu}_n' \boldsymbol{\Lambda}_n \boldsymbol{\mu}_n}{2}\right)$
 249 and $(g^{(j)}, f^{(j)})' | \dots \sim N(\boldsymbol{\mu}_n, \sigma^{2(j)} \boldsymbol{\Lambda}_n^{-1})$, with $\boldsymbol{\Lambda}_n = \mathbf{X}_{\mu^{(j)}}' \mathbf{X}_{\mu^{(j)}} + \boldsymbol{\Lambda}_0$ and $\boldsymbol{\mu}_n =$
 250 $\boldsymbol{\Lambda}_n^{-1} \left(\mathbf{X}_{\mu^{(j)}}' \mathbf{X}_{\mu^{(j)}} \left(\mathbf{X}_{\mu^{(j)}}' \mathbf{X}_{\mu^{(j)}} \right)^{-1} \mathbf{X}_{\mu^{(j)}}' \mathbf{y}_{\mu^{(j)}} + \boldsymbol{\Lambda}_0 \boldsymbol{\mu}_0 \right)$. The updates of $\mathbf{A}^{(j)}, \mathbf{b}^{(j)}$ and $\mathbf{Q}^{(j)}$ are
 251 similar, with diagonal assumption.