
Clustering Neural Populations by Poisson Dynamic Factor Model

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 [TODO]

2 Introduction

3 [TODO]

4 2 Method

5 2.1 Poisson Dynamic Factor Model

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

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

11 , 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
12 with Gaussian noise as

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

13 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
14 model can be equivalently written as regular Poisson factor model as

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

15 However, since the condition $T/N \rightarrow 0$ doesn't hold, the latent state $\mathbf{X}^{(j)}$ cannot be consistently
16 estimated, and assuming linear dynamics of $\mathbf{X}^{(j)}$ resolves the problem. Note that when $p > 1$, the

model is not unique, since $\widetilde{\mathbf{X}}^{(j)} = \mathbf{X}^{(j)}\mathbf{U}$ also satisfies the equation for any orthogonal matrix \mathbf{U} of order p . To ensure the model indefinability, we simply assume $\mathbf{A}^{(j)}$ and $\mathbf{Q}^{(j)}$ are both diagonal for convenience. See more detailed discussions of the constraints in discussion. Overall, denote the cluster-related parameters of cluster j as $\boldsymbol{\theta}^{(j)} = \{\boldsymbol{\mu}^{(j)}, \mathbf{X}^{(j)}, f^{(j)}, g^{(j)}, \Sigma^{(j)}, \{\mathbf{A}^{(j)}, \{\mathbf{b}^{(j)}, \{\mathbf{Q}^{(j)}\}$ and the spike counts of neuron i is generated by Poisson dynamic factor model (PDFM) as $\mathbf{Y}_i \sim PDFM(\boldsymbol{\theta}^{(z_i)})$, with the prior of $\boldsymbol{\theta}^{(j)}$ as \mathbf{H} . The priors for $\{f^{(j)}, g^{(j)}, \Sigma^{(j)}, \{\mathbf{A}^{(j)}, \{\mathbf{b}^{(j)}, \{\mathbf{Q}^{(j)}\}$ are regular normal and inverse-gamma distribution (or multivariate normal and inverse-Wishart when using other non-diagonal constraints).

The marginal likelihood of neuron i is

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

The marginal likelihood has no closed form and will be used for clustering. To help with fast clustering, instead of doing the Laplace approximation, we choose to make use of the Poisson-Gamma conjugacy [reference]. This leads to the closed form approximation. By the conditional independency assumption, $M_{\boldsymbol{\theta}^{(j)}}(\mathbf{y}_i) = \prod_{t=1}^T P(y_{it}|\boldsymbol{\theta}^{(j)})$. Since $\mathbf{c}_i \sim N(\mathbf{0}, \mathbf{I}_p)$, $\lambda_{it} = \exp(\mu_t^{(j)} + \mathbf{c}_i' \mathbf{x}_t^{(j)}) \sim \text{lognormal}(\mu_t^{(j)}, \mathbf{x}_t^{(j)} \mathbf{x}_t^{(j)})$. Approximate the lognormal distribution by Gamma distribution, s.t. $\text{lognormal}(\mu_t^{(j)}, \mathbf{x}_t^{(j)} \mathbf{x}_t^{(j)}) \approx \text{Gamma}(a_{it}, b_{it})$ with $a_{it} = (\mathbf{x}_t^{(j)} \mathbf{x}_t^{(j)})^{-1}$ and $b_{it} = \mathbf{x}_t^{(j)} \mathbf{x}_t^{(j)} \cdot e^{\mu_t^{(j)}}$. Then by Poisson-Gamma conjugacy,

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

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

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

2.2 Cluster by Mixture of Finite Mixtures Model

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

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

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

51 2.3 Inference

52 In this paper, we choose to do inference by MCMC. Because of the Poisson likelihood, the latent
53 state $\mathbf{X}^{(j)}$ has no closed full conditional distribution. We can sample the posterior by particle MCMC
54 directly, but this can be slow. However, due to the Markovian structure of the model, the conditional
55 log-posterior is concave and its Hessian is block-tridiagonal. Thus, we can do the global Laplace
56 approximation efficiently in $\mathcal{O}(T)$. The cluster index and number of cluster are sampled by the analog
57 of partition-based algorithm in DPM. See details of the MCMC in appendix.

58 In practice, using variational Bayes (VB) instead of MCMC may be more favorable. The PLDS can
59 be updated by variational EM. Using the stick-breaking representation of MFM model, we can do
60 VB easily similar to [reference David Blei & Michael I Jordan]. However, checking by the “gold
61 standard” MCMC before doing VB is always a good choice.

62 3 Simulations

63 3.1 Model Global Non-linearity by Clustering

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

70 We first simulated three neural populations, with 10 neurons in each. We set $p=1$ for each cluster and
71 the recording length is $T=1000$. After checking the trace plots up to 10,000 iterations, the convergence
72 achieved after several steps. The panel A and B in figure 1 show the posterior mean firing rate and
73 fitted latent state, averaging from iteration 500 to 1000. The latent state is transformed into the
74 commonly used PLDS model: $\log \lambda_i = \delta_i \mathbf{1}_T + \mathbf{G}^{(j)} \mathbf{d}_i$, where $\mathbf{d}_i \in \mathbb{R}^q$. In PLDS parametrization,
75 the latent state has one more dimension, i.e. $p = q + 1$. The PLDS parametrization also doesn't have
76 the unique solution. For comparison, we let $\mathbf{G}^{(j)}$ has mean zeros columns and is orthogonal.

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

85 3.1.1 Clustering

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

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

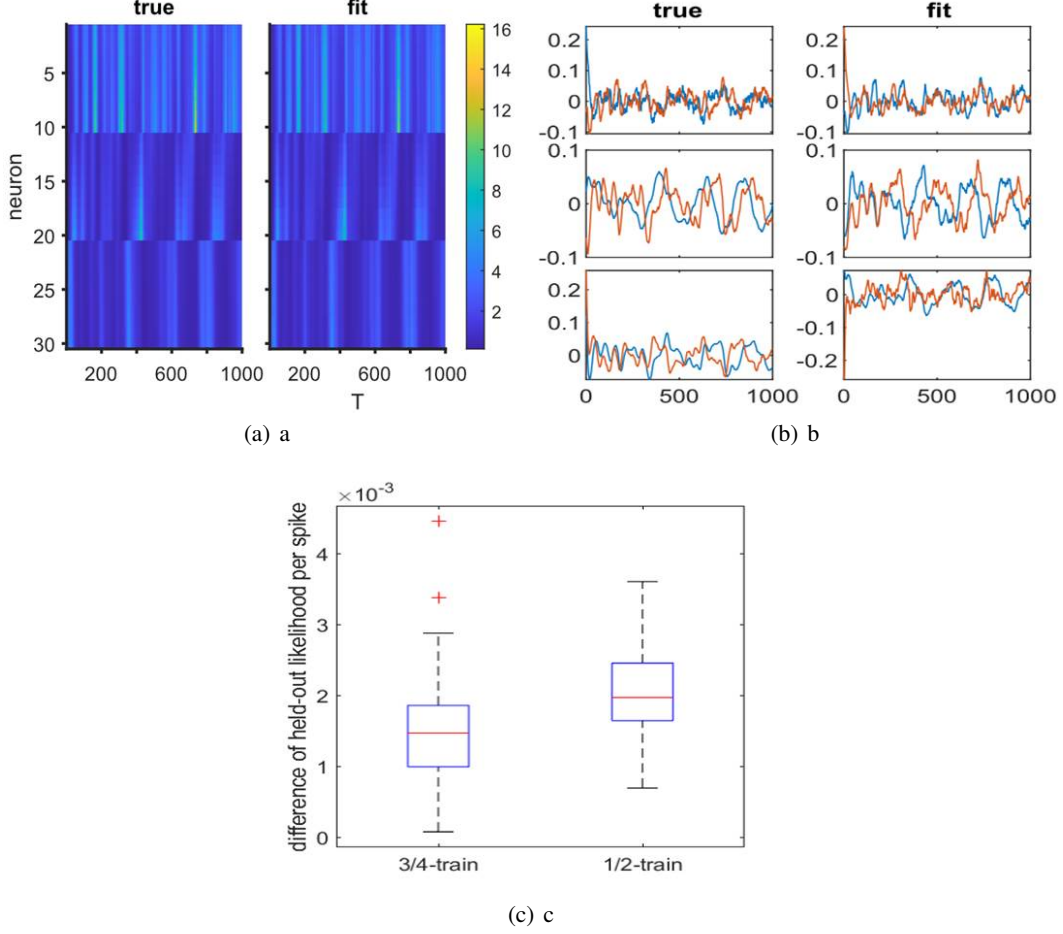


Figure 1: Model Global Non-linearity by Clustering

97 D in figure 2) of posterior, averaging from iteration 500 to 1000, shows that the model tend to split
 98 cluster 1 and 3 into two sub-clusters.

99 4 Application

100 4.1 Neuropixels data

101 **[TODO]**: brief introduction of Neuropixels dataset.

102 The Neuropixels dataset contains recording of neural activities in different brain regions, when doing
 103 the drift-grating experiment. Here, we only consider neurons with $\text{SNR} > 3$ and the spiking counts
 104 > 1000 . Then, we use the recording activity from Lateral posterior nucleus of the thalamus (LP, 20
 105 neurons), anteromedial visual area (VISam, 12 neurons) and ventral posteromedial nucleus of the
 106 thalamus (VPM, 14 neurons) during the spontaneous period. The bin size is 0.1s and truncate the
 107 data up to 1000 steps (100s recording) for convenience.

108 Panel A in figure 3 shows the spiking counts of these 46 neurons. Here we first fit model using
 109 all data, with $p = 1$ and $K \sim \text{Geometric}(0.3)$. For the number of clusters, panel B in figure 2
 110 shows the trace plot for the first 5000 iterations and panel C shows the histogram of iteration 2500 to
 111 5000. These plots show that these neurons are quite non-homogenous and they tend to form many
 112 sub-populations. The posterior similarity matrix, averaging from iteration 2500 to 5000 (panel D in

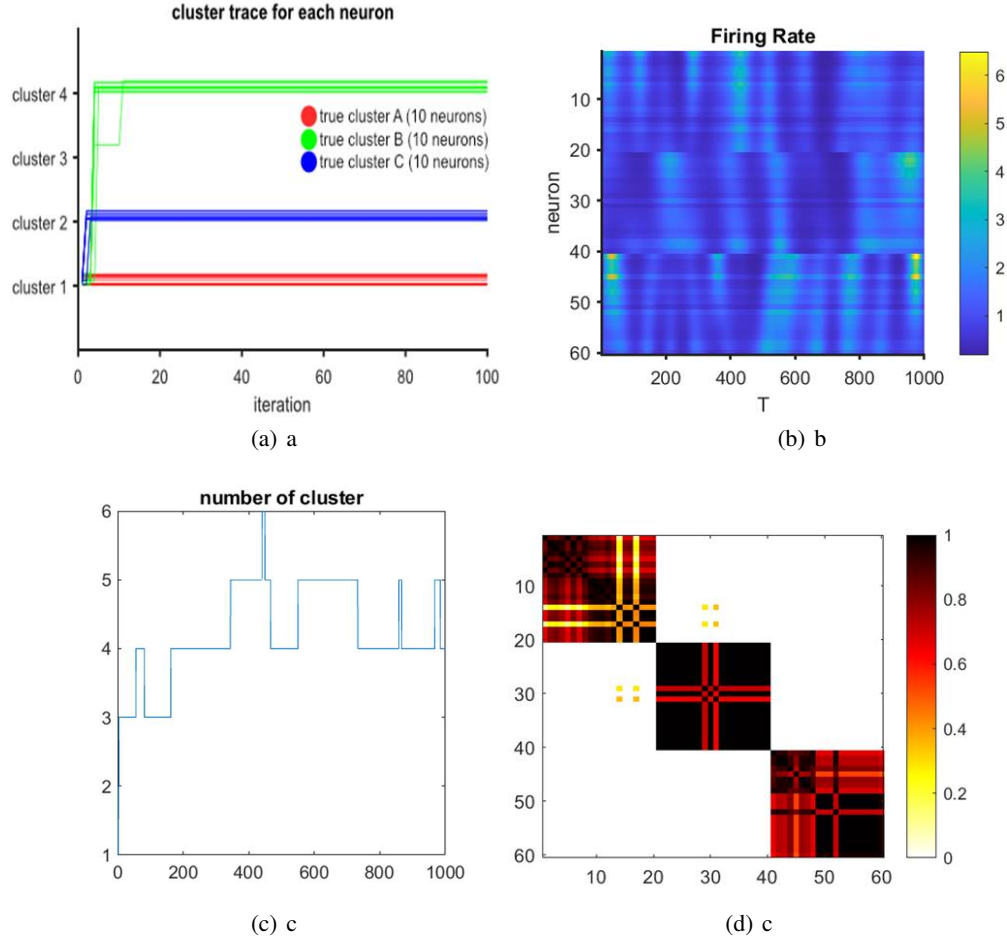


Figure 2: Clustering

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

[TODO]

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 A and 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.

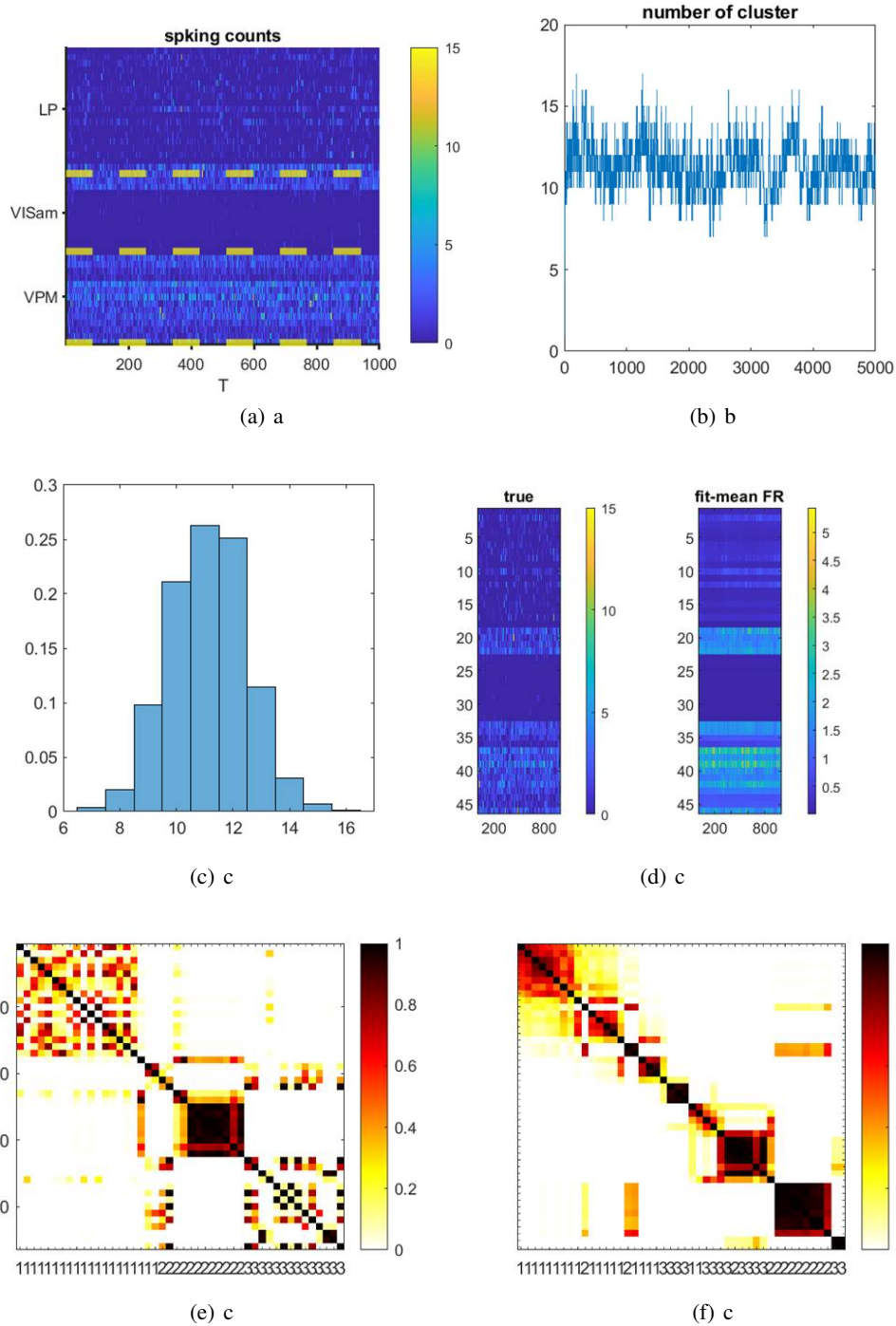


Figure 3: Neuropixel

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 $p(p-1)/2$ s 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: The first p rows of $\mathbf{X}^{(j)}$ is

136 lower triangular and the diagonal elements are positive. Since we assume $\mathbf{X}^{(j)}$ has linear progression,
137 using the diagonal $\mathbf{X}'^{(j)}\mathbf{X}^{(j)}$ constraint is more appropriate. This constraint can be easily achieved
138 when conducting deterministic optimization algorithms (e.g. the VB mentioned in “method-inference”
139 section). This can also be done in MCMC, by turning off the reflection of the projection (e.g. let the
140 diagonal elements in the projection matrix be positive). Using diagonal $\mathbf{X}'^{(j)}\mathbf{X}^{(j)}$ leads to similar
141 results.

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

147 **Acknowledgment**

148 **References**

- 149 [1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In
150 G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp.
151 609–616. Cambridge, MA: MIT Press.
- 152 [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the*
153 *GENeral NEural Simulation System*. New York: TELOS/Springer-Verlag.
- 154 [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent
155 synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.

156 **A Appendix**