
Appendix for Bayesian Clustering of Neural Spiking Activity Using a mixDPFA

A Appendix

A.1 MCMC updates

The posteriors are sampled by a Gibbs sampler. In each iteration, the sampling scheme has two main stages: 1) to sample the model parameters assuming known labels and 2) to sample the cluster indices given model parameters. Before moving into the second stage, the first stage is repeated several times. The sampling in the first stage is conducted without considering constraints for $\mu_t^{(j)}$ and $\mathbf{x}_t^{(j)}$ at first, and then we project the samples onto the constraint space for $\sum_{t=1}^T \mu_t^{(j)} = 0$ and $\sum_{t=1}^T \mathbf{x}_t^{(j)} = \mathbf{0}$ as used in (Sen et al., 2018).

Update $\mathbf{x}_t^{(j)}$ and $\mu_t^{(j)}$ The priors for initial population baseline and latent factor are:

$$\begin{aligned}\mu_1^{(j)} &\sim \mathcal{N}(0, 1), \\ \mathbf{x}_1^{(j)} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p).\end{aligned}$$

Assume there are $n_j = \#\{i : z_i = j\}$ neurons belong to the j -th cluster, and denote the spike counts and firing rates of these neurons at time t as $\tilde{\mathbf{y}}_t^{(j)} = \text{vec}(\{y_{it}|z_i = j\}) \in \mathbb{Z}_{\geq 0}^{n_j}$ and $\tilde{\boldsymbol{\lambda}}_t^{(j)} = \text{vec}(\{\lambda_{it}|z_i = j\}) \in \mathbb{R}^{n_j}$, where $\mathbb{Z}_{\geq 0}^{n_j}$ denotes a n_j -dimensional vector with each element being non-negative integers. The corresponding loading and baseline for these n_j neurons is $\mathbf{C}^{(j)} \in \mathbb{R}^{n_j \times p}$ and $\boldsymbol{\Delta}_j = \text{vec}(\{\delta_i|z_i = j\}) \in \mathbb{R}^{n_j}$, such that $\log \tilde{\boldsymbol{\lambda}}_t^{(j)} = \boldsymbol{\Delta}_j + \mu_t^{(j)} \mathbf{1}_{n_j} + \mathbf{C}^{(j)} \mathbf{x}_t^{(j)} = \boldsymbol{\Delta}_j + (\mathbf{1}_{n_j}, \mathbf{C}^{(j)}) \begin{pmatrix} \mu_t^{(j)} \\ \mathbf{x}_t^{(j)} \end{pmatrix}'$. Denote $\tilde{\mathbf{C}}^{(j)} = (\mathbf{1}_{n_j}, \mathbf{C}^{(j)})$, $\tilde{\boldsymbol{\lambda}}_t^{(j)} = \begin{pmatrix} \mu_t^{(j)} \\ \mathbf{x}_t^{(j)} \end{pmatrix}'$, $\tilde{\boldsymbol{\alpha}}^{(j)} = (\tilde{\boldsymbol{\alpha}}_1^{(j)}, \dots, \tilde{\boldsymbol{\alpha}}_T^{(j)})'$, $\tilde{\mathbf{A}}^{(j)} = \text{diag}(h^{(j)}, \mathbf{A}^{(j)})$, $\tilde{\mathbf{b}}^{(j)} = (g^{(j)}, \mathbf{b}'^{(j)})'$ and $\tilde{\mathbf{Q}}^{(j)} = \text{diag}(\sigma^{2(j)}, \mathbf{Q}^{(j)})$. The full conditional distribution $P(\tilde{\boldsymbol{\alpha}}^{(j)} | \dots) = P(\tilde{\boldsymbol{\alpha}}^{(j)} | \{\tilde{\mathbf{y}}_t^{(j)}\}_{t=1}^T, \boldsymbol{\Delta}_j, \tilde{\mathbf{C}}^{(j)}, \tilde{\mathbf{A}}^{(j)}, \tilde{\mathbf{b}}^{(j)}, \tilde{\mathbf{Q}}^{(j)})$ is approximated by a global Laplace approximation, i.e.,

$$\begin{aligned}P(\tilde{\boldsymbol{\alpha}}^{(j)} | \dots) &\approx \mathcal{N}_{(p+1)T}(\tilde{\boldsymbol{\alpha}}^{(j)} | \boldsymbol{\mu}_{\tilde{\boldsymbol{\alpha}}^{(j)}}, \boldsymbol{\Sigma}_{\tilde{\boldsymbol{\alpha}}^{(j)}}), \\ \boldsymbol{\mu}_{\tilde{\boldsymbol{\alpha}}^{(j)}} &= \arg \max_{\tilde{\boldsymbol{\alpha}}^{(j)}} P(\tilde{\boldsymbol{\alpha}}^{(j)} | \dots), \\ \boldsymbol{\Sigma}_{\tilde{\boldsymbol{\alpha}}^{(j)}} &= -(\nabla \nabla \log P(\tilde{\boldsymbol{\alpha}}^{(j)} | \dots))|_{\tilde{\boldsymbol{\alpha}}^{(j)} = \boldsymbol{\mu}_{\tilde{\boldsymbol{\alpha}}^{(j)}}})^{-1}.\end{aligned}$$

Then, taking the logarithm of the full conditional distribution, i.e., $\ell = \ell(\tilde{\boldsymbol{\alpha}}^{(j)}) = \log P(\tilde{\boldsymbol{\alpha}}^{(j)} | \dots)$, we have

$$\begin{aligned}\ell &= \text{const} + \sum_{t=1}^T \left(\tilde{\mathbf{y}}_t^{(j)'} (\boldsymbol{\Delta}_j + \tilde{\mathbf{C}}^{(j)} \tilde{\boldsymbol{\lambda}}_t^{(j)}) - \mathbf{1}_{n_j}' \tilde{\boldsymbol{\lambda}}_t^{(j)} \right) - \frac{1}{2} (\tilde{\boldsymbol{\alpha}}_1^{(j)} - \tilde{\boldsymbol{\alpha}}_0^{(j)}) [\tilde{\mathbf{Q}}_0^{(j)}]^{-1} (\tilde{\boldsymbol{\alpha}}_1^{(j)} - \tilde{\boldsymbol{\alpha}}_0^{(j)}) \\ &\quad - \sum_{t=2}^T \frac{1}{2} (\tilde{\boldsymbol{\alpha}}_t^{(j)} - \tilde{\mathbf{A}}^{(j)} \tilde{\boldsymbol{\alpha}}_{t-1}^{(j)} - \tilde{\mathbf{b}}^{(j)}) [\tilde{\mathbf{Q}}^{(j)}]^{-1} (\tilde{\boldsymbol{\alpha}}_t^{(j)} - \tilde{\mathbf{A}}^{(j)} \tilde{\boldsymbol{\alpha}}_{t-1}^{(j)} - \tilde{\mathbf{b}}^{(j)}).\end{aligned}$$

Here, $\ell(\tilde{\mathbf{x}}^{(j)})$ is concave and unimodal. By the Markovian assumption for latent state vectors, the Hessian matrix is tri-block diagonal. We can thus compute Newton updates to get $\mu_{\tilde{\mathbf{x}}^{(j)}}$ for the Laplace approximation in $\mathcal{O}(T)$ (Paninski et al., 2010), similar to the E-step for the PLDS (Macke et al., 2011).

The gradient $\nabla\ell$ and the Hessian $\nabla\nabla\ell$ are provided as follows. For $t = 2, \dots, T-1$, the gradient of $\ell(\tilde{\mathbf{x}}^{(j)})$ is:

$$\begin{aligned}\nabla\ell &= \left[\left(\frac{\partial\ell}{\partial\tilde{\mathbf{x}}_1^{(j)}} \right)', \dots, \left(\frac{\partial\ell}{\partial\tilde{\mathbf{x}}_T^{(j)}} \right)' \right]', \\ \frac{\partial\ell}{\partial\tilde{\mathbf{x}}_1^{(j)}} &= \tilde{\mathbf{C}}'^{(j)}(\tilde{\mathbf{y}}_1^{(j)} - \tilde{\boldsymbol{\lambda}}_1) - [\tilde{\mathbf{Q}}_0^{(j)}]^{-1}(\tilde{\mathbf{x}}_1^{(j)} - \tilde{\mathbf{x}}_0^{(j)}) + \tilde{\mathbf{A}}'^{(j)}[\tilde{\mathbf{Q}}^{(j)}]^{-1}(\tilde{\mathbf{x}}_2^{(j)} - \tilde{\mathbf{A}}^{(j)}\tilde{\mathbf{x}}_1^{(j)} - \tilde{\mathbf{b}}^{(j)}), \\ \frac{\partial\ell}{\partial\tilde{\mathbf{x}}_t^{(j)}} &= \tilde{\mathbf{C}}'^{(j)}(\tilde{\mathbf{y}}_t^{(j)} - \tilde{\boldsymbol{\lambda}}_t) - [\tilde{\mathbf{Q}}^{(j)}]^{-1}(\tilde{\mathbf{x}}_t^{(j)} - \tilde{\mathbf{A}}^{(j)}\tilde{\mathbf{x}}_{t-1}^{(j)} - \tilde{\mathbf{b}}^{(j)}) \\ &\quad + \tilde{\mathbf{A}}'^{(j)}[\tilde{\mathbf{Q}}^{(j)}]^{-1}(\tilde{\mathbf{x}}_{t+1}^{(j)} - \tilde{\mathbf{A}}^{(j)}\tilde{\mathbf{x}}_t^{(j)} - \tilde{\mathbf{b}}^{(j)}), \\ \frac{\partial\ell}{\partial\tilde{\mathbf{x}}_T^{(j)}} &= \tilde{\mathbf{C}}'^{(j)}(\tilde{\mathbf{y}}_T^{(j)} - \tilde{\boldsymbol{\lambda}}_T) - [\tilde{\mathbf{Q}}^{(j)}]^{-1}(\tilde{\mathbf{x}}_T^{(j)} - \tilde{\mathbf{A}}^{(j)}\tilde{\mathbf{x}}_{T-1}^{(j)} - \tilde{\mathbf{b}}^{(j)}).\end{aligned}$$

And the Hessian matrix is:

$$\begin{aligned}\nabla\nabla\ell &= \begin{pmatrix} \frac{\partial^2\ell}{\partial\tilde{\mathbf{x}}_1^{(j)}\partial\tilde{\mathbf{x}}_1^{(j)}} & \tilde{\mathbf{A}}'^{(j)}[\tilde{\mathbf{Q}}^{(j)}]^{-1} & \mathbf{0} & \dots & \mathbf{0} \\ [\tilde{\mathbf{Q}}^{(j)}]^{-1}\tilde{\mathbf{A}}^{(j)} & \frac{\partial^2\ell}{\partial\tilde{\mathbf{x}}_2^{(j)}\partial\tilde{\mathbf{x}}_2^{(j)}} & \tilde{\mathbf{A}}'^{(j)}[\tilde{\mathbf{Q}}^{(j)}]^{-1} & \dots & \vdots \\ \mathbf{0} & [\tilde{\mathbf{Q}}^{(j)}]^{-1}\tilde{\mathbf{A}}^{(j)} & \frac{\partial^2\ell}{\partial\tilde{\mathbf{x}}_3^{(j)}\partial\tilde{\mathbf{x}}_3^{(j)}} & \dots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \dots & \dots & \dots & \frac{\partial^2\ell}{\partial\tilde{\mathbf{x}}_T^{(j)}\partial\tilde{\mathbf{x}}_T^{(j)}} \end{pmatrix}, \\ \frac{\partial^2\ell}{\partial\tilde{\mathbf{x}}_1^{(j)}\partial\tilde{\mathbf{x}}_1^{(j)}} &= -\tilde{\mathbf{C}}'^{(j)}\text{diag}(\tilde{\boldsymbol{\lambda}}_1)\tilde{\mathbf{C}}^{(j)} - [\tilde{\mathbf{Q}}_0^{(j)}]^{-1} - \tilde{\mathbf{A}}'^{(j)}[\tilde{\mathbf{Q}}^{(j)}]^{-1}\tilde{\mathbf{A}}^{(j)}, \\ \frac{\partial^2\ell}{\partial\tilde{\mathbf{x}}_t^{(j)}\partial\tilde{\mathbf{x}}_t^{(j)}} &= -\tilde{\mathbf{C}}'^{(j)}\text{diag}(\tilde{\boldsymbol{\lambda}}_t)\tilde{\mathbf{C}}^{(j)} - [\tilde{\mathbf{Q}}^{(j)}]^{-1} - \tilde{\mathbf{A}}'^{(j)}[\tilde{\mathbf{Q}}^{(j)}]^{-1}\tilde{\mathbf{A}}^{(j)}, \\ \frac{\partial^2\ell}{\partial\tilde{\mathbf{x}}_T^{(j)}\partial\tilde{\mathbf{x}}_T^{(j)}} &= -\tilde{\mathbf{C}}'^{(j)}\text{diag}(\tilde{\boldsymbol{\lambda}}_T)\tilde{\mathbf{C}}^{(j)} - [\tilde{\mathbf{Q}}^{(j)}]^{-1}.\end{aligned}$$

Although we can conduct the Newton update efficiently in $\mathcal{O}(T)$, bad initial values may slow down the convergence. To facilitate convergence, we initialize the Newton update with a smoothing estimate by local Gaussian approximation. The forward filtering for a dynamic Poisson model has been previously described (Eden et al., 2004), and we use an additional backward pass to smooth (Rauch et al., 1965).

Let $\tilde{\mathbf{x}}_{t|t-1}^{(j)} = E(\tilde{\mathbf{x}}_t^{(j)}|\tilde{\mathbf{y}}_1^{(j)}, \dots, \tilde{\mathbf{y}}_{t-1}^{(j)})$ and $\mathbf{V}_{t|t-1}^{(j)} = \text{Var}(\tilde{\mathbf{x}}_t^{(j)}|\tilde{\mathbf{y}}_1^{(j)}, \dots, \tilde{\mathbf{y}}_{t-1}^{(j)})$ be the mean and variance for the one-step prediction density, where $\tilde{\mathbf{x}}_{t|t-1}^{(j)} = \tilde{\mathbf{A}}^{(j)}\tilde{\mathbf{x}}_{t-1|t-1}^{(j)} + \tilde{\mathbf{b}}^{(j)}$ and $\mathbf{V}_{t|t-1}^{(j)} = \tilde{\mathbf{A}}^{(j)}\mathbf{V}_{t-1|t-1}^{(j)}\tilde{\mathbf{A}}'^{(j)} + \tilde{\mathbf{Q}}^{(j)}$. Then, after we observe the data at time t , we can do a forward filtering step for the mean $\tilde{\mathbf{x}}_{t|t}^{(j)} = E(\tilde{\mathbf{x}}_t^{(j)}|\tilde{\mathbf{y}}_1^{(j)}, \dots, \tilde{\mathbf{y}}_t^{(j)})$ and the variance $\mathbf{V}_{t|t}^{(j)} = \text{Var}(\tilde{\mathbf{x}}_t^{(j)}|\tilde{\mathbf{y}}_1^{(j)}, \dots, \tilde{\mathbf{y}}_t^{(j)})$, which are given by

$$\begin{aligned}\tilde{\mathbf{x}}_{t|t}^{(j)} &= \tilde{\mathbf{x}}_{t|t-1}^{(j)} + \mathbf{V}_{t|t-1}^{(j)}[\tilde{\mathbf{C}}'^{(j)}(\tilde{\mathbf{y}}_t^{(j)} - \tilde{\boldsymbol{\lambda}}_t)]_{\tilde{\mathbf{x}}_t^{(j)} = \tilde{\mathbf{x}}_{t|t-1}^{(j)}}, \\ [\mathbf{V}_{t|t}^{(j)}]^{-1} &= [\mathbf{V}_{t|t-1}^{(j)}]^{-1} + [\tilde{\mathbf{C}}'^{(j)}\text{diag}(\tilde{\boldsymbol{\lambda}}_t)\tilde{\mathbf{C}}^{(j)} - [\tilde{\mathbf{Q}}^{(j)}]^{-1}]_{\tilde{\mathbf{x}}_t^{(j)} = \tilde{\mathbf{x}}_{t|t-1}^{(j)}}.\end{aligned}$$

Derivation of the filtering estimates can be found in (Eden et al., 2004), and we can further get the smoothing estimates directly by standard Rauch-Tung-Striebel smoother (Rauch et al., 1965).

The smoother estimates $\tilde{\mathbf{x}}_{t|T}^{(j)}$ and $\mathbf{V}_{t|T}^{(j)}$ are updated as follows:

$$\begin{aligned}\tilde{\mathbf{x}}_{t-1|T}^{(j)} &= \tilde{\mathbf{x}}_{t-1|t-1}^{(j)} + \mathbf{J}_{t-1}(\tilde{\mathbf{x}}_{t|T}^{(j)} - \tilde{\mathbf{x}}_{t|t-1}^{(j)}), \\ \mathbf{V}_{t-1|T}^{(j)} &= \mathbf{V}_{t-1|t-1}^{(j)} + \mathbf{J}_{t-1}(\mathbf{V}_{t|T}^{(j)} - \mathbf{V}_{t|t-1}^{(j)})\mathbf{J}'_{t-1},\end{aligned}$$

where $\mathbf{J}_{t-1} = \mathbf{V}_{t-1|t-1}^{(j)}\tilde{\mathbf{A}}'^{(j)}[\mathbf{V}_{t|t-1}^{(j)}]^{-1}$.

Update δ_i and \mathbf{c}_i We specify the prior for neuron-specific baseline δ_i as $\delta_i \sim \mathcal{N}(0, 1)$ and we have assumed the loading $\mathbf{c}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$. Then, from the matrix representation of DPFA in Equation (1), i.e., $\log \lambda_i = \delta_i \mathbf{1}_T + \boldsymbol{\mu}^{(j)} + \mathbf{X}^{(j)}\mathbf{c}_i$, it is easy to see that given $\boldsymbol{\mu}^{(j)}$ and $\mathbf{X}^{(j)}$ are known, the update of δ_i and \mathbf{c}_i is just a regular Bayesian Poisson regression problem. Thus, we can sample the full conditional distribution of δ_i and \mathbf{c}_i by a Hamiltonian Monte Carlo (HMC, Duane et al. 1987) within the Gibbs sampler.

Update parameters of latent state The parameters for linear dynamics are $h^{(j)}$, $g^{(j)}$, $\sigma^{2(j)}$, $\mathbf{A}^{(j)}$, $\mathbf{b}^{(j)}$ and $\mathbf{Q}^{(j)}$. To make the model identifiable, we simply assume $\mathbf{A}^{(j)} = \text{diag}(a_1^{(j)}, \dots, a_p^{(j)})$ and $\mathbf{Q}^{(j)} = \text{diag}(q_1^{(j)}, \dots, q_p^{(j)})$. Therefore, we can update $\mathbf{A}^{(j)}$, $\mathbf{b}^{(j)}$ and $\mathbf{Q}^{(j)}$ for each diagonal element separately, as the update in $h^{(j)}$, $g^{(j)}$ and $\sigma^{2(j)}$. Here, we update $h^{(j)}$, $g^{(j)}$ and $\sigma^{2(j)}$ as follows.

First, we specify the priors for $\sigma^{2(j)}$ following $IG(\nu_0/2, \nu_0\sigma_0^2/2)$ and $(g^{(j)}, h^{(j)})' \sim \mathcal{N}(\boldsymbol{\tau}_0, \sigma^{2(j)}\boldsymbol{\Lambda}_0^{-1})$, with $\nu_0 = 1$, $\sigma_0 = 0.01$, $\boldsymbol{\tau}_0 = (0, 1)'$ and $\boldsymbol{\Lambda}_0 = \mathbf{I}_2$. Here, the "IG" denotes the inverse-gamma distribution.

Denote $\boldsymbol{\mu}_{2:T}^{(j)} = (\mu_2^{(j)}, \dots, \mu_T^{(j)})'$ and $\tilde{\boldsymbol{\mu}}_{1:(T-1)}^{(j)} = (\mathbf{1}_{T-1}, \boldsymbol{\mu}_{1:(T-1)}^{(j)})$, with $\boldsymbol{\mu}_{1:(T-1)}^{(j)} = (\mu_1^{(j)}, \dots, \mu_{T-1}^{(j)})'$. The full conditional distributions for $\sigma^{2(j)}$ and $(g^{(j)}, h^{(j)})'$ are:

$$\begin{aligned}\sigma^{2(j)} | \{\mu_t^{(j)}\}_{t=1}^T &\sim IG\left(\frac{\nu_0 + T - 1}{2}, \frac{\nu_0\sigma_0^2 + \boldsymbol{\mu}_{2:T}'^{(j)}\boldsymbol{\mu}_{2:T}^{(j)} + \boldsymbol{\tau}_0'\boldsymbol{\Lambda}_0\boldsymbol{\tau}_0 - \boldsymbol{\tau}_n'\boldsymbol{\Lambda}_n\boldsymbol{\tau}_n}{2}\right), \\ (g^{(j)}, h^{(j)})' | \{\mu_t^{(j)}\}_{t=1}^T &\sim \mathcal{N}(\boldsymbol{\tau}_n, \sigma^{2(j)}\boldsymbol{\Lambda}_n^{-1}),\end{aligned}$$

with $\boldsymbol{\Lambda}_n = \tilde{\boldsymbol{\mu}}_{1:(T-1)}^{(j)}\tilde{\boldsymbol{\mu}}_{1:(T-1)}^{(j)} + \boldsymbol{\Lambda}_0$, and $\boldsymbol{\tau}_n = \boldsymbol{\Lambda}_n^{-1}(\tilde{\boldsymbol{\mu}}_{1:(T-1)}^{(j)}\boldsymbol{\mu}_{2:T}^{(j)} + \boldsymbol{\Lambda}_0\boldsymbol{\tau}_0)$.

For completeness, we also provide the update of latent state parameters when using the more parsimonious constraint, i.e. diagonal $\mathbf{X}'^{(j)}\mathbf{X}^{(j)}$. According to results from previous research (Krzanowski and Marriott, 1994b,a; Fokoué and Titterton, 2003), this constraint is one of the most parsimonious one, which only put constraints on $p(p-1)/2$ parameters. The constraint can also be implemented in MCMC by "unconstrained sampling-projection" procedure (Sen et al., 2018).

Without constraints, the sampling of $h^{(j)}$, $g^{(j)}$ and $\sigma^{2(j)}$ is the same as shown previously. The update of $\mathbf{A}^{(j)}$, $\mathbf{b}^{(j)}$ and $\mathbf{Q}^{(j)}$ is the standard multivariate Bayesian linear regression. Denote $\mathbf{X}_{2:T}^{(j)} = (\mathbf{x}_2^{(j)}, \dots, \mathbf{x}_T^{(j)})'$ and $\tilde{\mathbf{X}}_{1:(T-1)}^{(j)} = (\mathbf{1}_{T-1}, \mathbf{X}_{1:(T-1)}^{(j)})$, with $\mathbf{X}_{1:(T-1)}^{(j)} = (\mathbf{x}_1^{(j)}, \dots, \mathbf{x}_{T-1}^{(j)})'$. Let us use the conjugate priors as following for $\mathbf{Q}^{(j)}$, $\mathbf{b}^{(j)}$ and $\mathbf{A}^{(j)}$:

$$\begin{aligned}\mathbf{Q}^{(j)} &\sim \mathcal{W}^{-1}(\boldsymbol{\Psi}_0, \gamma_0), \\ \text{vec}((\mathbf{b}^{(j)}, \mathbf{A}^{(j)})) &\sim N(\text{vec}(\mathbf{T}_0), \mathbf{Q}^{(j)} \otimes \boldsymbol{\Gamma}_0^{-1}).\end{aligned}$$

Here, the \mathcal{W}^{-1} denotes the inverse-Wishart distribution. We can set the priors as $\boldsymbol{\Psi}_0 = 0.01\mathbf{I}_p$, $\gamma_0 = p+2$ and $\mathbf{T}_0 = (\mathbf{0}_p, \mathbf{I}_p)'$. Then, the full conditional distributions for $\mathbf{Q}^{(j)}$ and $\text{vec}((\mathbf{b}^{(j)}, \mathbf{A}^{(j)}))'$ are:

$$\begin{aligned}\mathbf{Q}^{(j)} | \mathbf{X}^{(j)} &\sim \mathcal{W}^{-1}(\boldsymbol{\Psi}_n, \gamma_n), \\ \text{vec}((\mathbf{b}^{(j)}, \mathbf{A}^{(j)}))' | \mathbf{X}^{(j)} &\sim N(\text{vec}(\mathbf{T}_n), \mathbf{Q}^{(j)} \otimes \boldsymbol{\Gamma}_n^{-1}),\end{aligned}$$

with

$$\begin{aligned}
\boldsymbol{\Psi}_n &= \boldsymbol{\Psi}_0 + (\mathbf{X}_{2:T}^{(j)} - \widetilde{\mathbf{X}}_{1:(T-1)}^{(j)} \mathbf{T}_n)' (\mathbf{X}_{2:T}^{(j)} - \widetilde{\mathbf{X}}_{1:(T-1)}^{(j)} \mathbf{T}_n), \\
&\quad + (\mathbf{T}_n - \mathbf{T}_0)' \boldsymbol{\Gamma}_0 (\mathbf{T}_n - \mathbf{T}_0), \\
\gamma_n &= \gamma_0 + T - 1, \\
\mathbf{T}_n &= \boldsymbol{\Gamma}_n^{-1} (\widetilde{\mathbf{X}}_{1:(T-1)}^{(j)} \mathbf{X}_{2:T}^{(j)} + \boldsymbol{\Gamma}_0 \mathbf{T}_0), \\
\boldsymbol{\Gamma}_n &= \widetilde{\mathbf{X}}_{1:(T-1)}^{(j)} \widetilde{\mathbf{X}}_{1:(T-1)}^{(j)} + \boldsymbol{\Gamma}_0.
\end{aligned}$$

Before updating the cluster indices z_i , the sampling of $\{\mathbf{x}_t^{(j)}, \mu_t^{(j)}, \delta_i, \mathbf{c}_i, \boldsymbol{\theta}^{(j)}\}$ is repeated several times (5 times in both simulation and application for this paper). To save time, we can further allow the repetitions to be pre-stopped, when the training log-likelihood converges roughly.

Update z_i To update the cluster assignments for each neuron i , we use a partition based algorithm for MFM, similarly as described in Miller and Harrison 2018.

Let \mathcal{C} denote a partition of neurons, and $\mathcal{C} \setminus i$ denote the partition obtained by removing neuron i from \mathcal{C} .

1. Initialize \mathcal{C} and $\{\boldsymbol{\theta}^{(c)} : c \in \mathcal{C}\}$ (e.g. one cluster).
2. Repeat the following steps G times to obtain G samples. For $i = 1, \dots, N$: remove neuron i from \mathcal{C} and place it:
 - (a) in $c \in \mathcal{C} \setminus i$ with probability $\propto (|c| + \gamma) M_{\boldsymbol{\theta}^{(c)}}(\mathbf{y}_i)$, where γ is defined in the MFM model in the main text (Equation 2) and $M_{\boldsymbol{\theta}^{(c)}}(\mathbf{y}_i)$ denotes the marginal likelihood of neuron i in cluster c , when integrating the loading \mathbf{c}_i out (Equation (3)). The marginal likelihood is approximated by using Poisson-Gamma conjugacy.
 - (b) in a new cluster c^* with probability $\propto \gamma \frac{V_n(t+1)}{V_n(t)} M_{\boldsymbol{\theta}^{(c^*)}}(\mathbf{y}_i)$, where t is the number of partitions obtained by removing the neuron i and $V_n(t) = \sum_{j=1}^{\infty} \frac{j^{(t)}}{(\gamma j)^{(n)}} f_k(j)$, with $x^{(m)} = x(x+1) \cdots (x+m-1)$, $x_{(m)} = x(x-1) \cdots (x-m+1)$, $x^{(0)} = 1$ and $x_{(0)} = 1$.

The update is an adaptation of partition-based algorithm for DPM (Neal, 2000), but with two substitutions: 1) replace $|c_i|$ by $|c_i| + \gamma$ and 2) replace α by $\gamma V_n(t+1)/V_n(t)$. See more details and discussions in (Miller and Harrison, 2018). When evaluating the likelihood, we marginalize the cluster-independent loading \mathbf{c}_i out. This is necessary for the high dimensional situation, otherwise the chain will stop moving.

One issue with incremental Gibbs samplers such as Algorithm 3 and 8 in Neal (2000), when applied to DPM, is that mixing can be somewhat slow. To further improve the mixing, we may intersperse the "split-merge" Metropolis-Hasting updates (Jain and Neal, 2007, 2004) between Gibbs sweeps, as in (Miller and Harrison, 2018).

A.2 Sample Latent Vectors Using Pólya-Gamma Augmentation and Metropolis-Hastings Algorithm

In this section, we provide details of sampling algorithms to draw the latent states $\mathbf{X}^{(j)}$ and population baseline $\boldsymbol{\mu}^{(j)}$ from the exact posterior of the model. In order to explain the algorithm much clearer, we just focus on a given cluster index, and thus in this section we suspend the superscript (j) as the cluster index in our notations. We present the sampling idea with two major parts as described below.

Pólya-Gamma Augmentation Although the mixDPFA model does not directly follow the PG augmentation scheme (Polson et al., 2013), we can approximate the Poisson distribution by a negative binomial (NB) distribution, i.e., $\lim_{r \rightarrow \infty} \text{NB}(r, \sigma(\psi - \log r)) = \text{Poisson}(e^\psi)$, where $\sigma(\psi) = e^\psi / (1 + e^\psi)$ and $\text{NB}(r, p)$ denotes the NB distribution with $rp/(1-p)$ as its expectation. We approximate the proposed DPFA model using a NB, then we can instead sample a mixture of NB factor analyzers using the PG scheme (Windle et al., 2013). Further, we use the forward-filtering-backward-sampling (FFBS) algorithm (Carter and Kohn, 1994; Frühwirth-Schnatter, 1994) to update

the latent states \mathbf{X} and population baseline $\boldsymbol{\mu}$. These two steps are summarized in the step 1 and step 2 of Algorithm 1 for updating $\widetilde{\mathbf{X}} = (\boldsymbol{\mu}, \mathbf{X})$.

Metropolis-Hastings Step Next, we use the samples of $\widetilde{\mathbf{X}}$ yielded from FFBS algorithm as a proposal, where we employ a Metropolis-Hastings (MH) step to reject or accept the proposal. In this step, the dispersion parameter r in NB distribution becomes a tuning parameter, to balance acceptance rate and autocorrelation in MH. When r is large, the approximation to Poisson observation is accurate and the MH performs similar to the Gibbs sampler. Here, we allow neurons at different time points to have unique tuning parameters r_{it} . The MH step is summarized in step 3 of Algorithm 1.

We have further implemented Algorithm 1 to fit DPFA model for cluster 1 data in Fig. 2 and the results have been presented in Fig. 1. We set the tuning parameters for each neuron at all time points as $r_{it} = 10$, to tune the acceptance rate around 0.4. When implementing the PG-MH algorithm in DPFA, the ω_{it} (i.e., PG random variable) was sampled using `pgdraw` function in R package `pgdraw`, which is called from our MATLAB code. Although we can program Algorithm 1 without calling `pgdraw` function from R to save some time, to sample the PG random variable is often much computationally expensive in comparison to sample from the Gaussian random variable. When fitting the simulated data with $N = 5$, $T = 1000$ and $p = 2$ (cluster 1 data in Fig. 2) using a 3.40 GHz processor with 16 GB of RAM, it takes 131.64s for PG-MH to draw 1000 posterior samples, while the proposed approximation method in the main context of the paper takes 83.86s. Although the proposed method takes the approximation to the full conditional of the latent state \mathbf{X} and population baseline $\boldsymbol{\mu}$, the results are similar as to sample directly from the exact full conditional. Table 1 shows the similarities, defined by cosine function, of fitted $\boldsymbol{\mu}$ and $\mathbf{X} = (\mathbf{X}_{*1}, \mathbf{X}_{*2})$ between the PG-MH and the Gaussian approximation for full conditional, where \mathbf{X}_{*l} denotes the l -th latent vectors (i.e. the l -th column of \mathbf{X}). Here, $\hat{\boldsymbol{\xi}}_{\text{PGMH}}$ denotes the average from iteration 5000 to 10,000 for chain 1 of PG-MH (Fig. 1), $\hat{\boldsymbol{\xi}}_{\mathcal{N}}$ denotes the average from iteration 5000 to 10,000 for the chain in Fig. 2 of the Gaussian approximation, and $\boldsymbol{\xi}_{\text{true}}$ denotes the ground truth. In each column, $\boldsymbol{\xi}$ is replaced by $\boldsymbol{\mu}$, \mathbf{X}_{*1} and \mathbf{X}_{*2} respectively. The posterior means for these two method are close to each other, i.e. $\cos(\hat{\boldsymbol{\xi}}_{\text{PGMH}}, \hat{\boldsymbol{\xi}}_{\mathcal{N}}) \approx 1$.

A.3 Supplementary Results for Neuropixels Application

This section shows supplementary results when applying the proposed model to the Neuropixels dataset (4 Multi-region neural spike recordings). Here, we show 1) results sorted by maximum a posteriori probability (MAP) estimates and 2) clustering results in another independent chain (Fig. 2).

Algorithm 1: Pólya-Gamma-Metropolis-Hastings Algorithm (PG-MH) for Poisson Dynamic Model

Recall the DPFA model using the notations defined in A.1:

$$y_{it} \sim Poi(\lambda_{it}),$$

$$\log \lambda_{it} = \delta_i + \tilde{\mathbf{c}}_i' \tilde{\mathbf{x}}_t,$$

for $i = 1, \dots, n$ and $t = 1, \dots, T$. Here, $\tilde{\mathbf{c}}_i = (1, \mathbf{c}'_i)'$ and $\tilde{\mathbf{x}}_t = (\mu_t, \mathbf{x}'_t)'$. $\tilde{\mathbf{x}}_t$ follows linear dynamics $\tilde{\mathbf{x}}_{t+1} | \tilde{\mathbf{x}}_t \sim \mathcal{N}(\tilde{\mathbf{A}}\tilde{\mathbf{x}}_t + \tilde{\mathbf{b}}, \tilde{\mathbf{Q}})$. Denote the prior as $\tilde{\mathbf{x}}_1 \sim \mathcal{N}(\mathbf{m}_0, \mathbf{V}_0)$. Given the sample from the $(G-1)$ -th iteration $\tilde{\mathbf{x}}_t^{(G-1)}$ and $\mathbf{U} = \{\tilde{\mathbf{c}}_i, \delta_i, \tilde{\mathbf{A}}, \tilde{\mathbf{b}}, \tilde{\mathbf{Q}}\}$.

1. sample ω_{it} from PG distribution and calculate \hat{y}_{it} , which follows $\mathcal{N}(\tilde{\mathbf{c}}_i' \tilde{\mathbf{x}}_t^{(G-1)}, \omega_{it}^{-1})$

for $t = 1, \dots, T$ **do**

for $i = 1, \dots, n$ **do**

$$\begin{aligned} & \text{sample } \omega_{it} \sim PPG(r_{it} + y_{it}, \delta_i + \tilde{\mathbf{c}}_i' \tilde{\mathbf{x}}_t^{(G-1)} - \log r_{it}) \\ & \kappa_{it} = (y_{it} - r_{it})/2 + \omega_{it}(\log r_{it} - \delta_i) \\ & \hat{y}_{it} = \omega_{it}^{-1} \kappa_{it} \end{aligned}$$

end

end

2. Forward-filtering-backward-sampling (FFBS) for $\tilde{\mathbf{X}}$

Denote $\hat{\mathbf{y}}_t = (\hat{y}_{1t}, \dots, \hat{y}_{Nt})'$, $\mathbf{\Omega}_t = \text{Diag}([\omega_{1t}, \dots, \omega_{Nt}])$ and $\tilde{\mathbf{C}} = (\tilde{\mathbf{c}}_1, \dots, \tilde{\mathbf{c}}_N)'$

for $t = 1, \dots, T$ **do**

$$\begin{aligned} & \mathbf{m}_{t|t-1} = \tilde{\mathbf{A}}\mathbf{m}_{t-1} + \tilde{\mathbf{b}} \\ & \mathbf{V}_{t|t-1} = \tilde{\mathbf{A}}\mathbf{V}_{t-1}\tilde{\mathbf{A}}' + \tilde{\mathbf{Q}} \\ & \mathbf{K}_t = \mathbf{V}_{t|t-1}\tilde{\mathbf{C}}'(\tilde{\mathbf{C}}\mathbf{V}_{t|t-1}\tilde{\mathbf{C}}' + \mathbf{\Omega}_t^{-1})^{-1} \\ & \mathbf{m}_t = \mathbf{m}_{t|t-1} + \mathbf{K}_t(\hat{\mathbf{y}}_t - \tilde{\mathbf{C}}\mathbf{m}_{t|t-1}) \\ & \mathbf{V}_t = (\mathbf{I} - \mathbf{K}_t\tilde{\mathbf{C}})\mathbf{V}_{t|t-1} \end{aligned}$$

end

sample $\tilde{\mathbf{x}}_T^* \sim \mathcal{N}(\mathbf{m}_T, \mathbf{V}_T)$

for $t = T-1, \dots, 1$ **do**

$$\begin{aligned} & \mathbf{J}_t = \mathbf{V}_t\tilde{\mathbf{A}}'(\tilde{\mathbf{A}}\mathbf{V}_t\tilde{\mathbf{A}}' + \tilde{\mathbf{Q}})^{-1} \\ & \mathbf{m}_t^* = \mathbf{m}_t + \mathbf{J}_t(\tilde{\mathbf{x}}_{t+1}^* - \tilde{\mathbf{A}}\mathbf{m}_t - \tilde{\mathbf{b}}) \\ & \mathbf{V}_t^* = (\mathbf{I} - \mathbf{J}_t\tilde{\mathbf{A}})\mathbf{V}_t \\ & \text{sample } \tilde{\mathbf{x}}_t^* \sim \mathcal{N}(\mathbf{m}_t^*, \mathbf{V}_t^*) \end{aligned}$$

end

3. Accept or reject the proposal $\tilde{\mathbf{X}}^*$

compute the acceptance ratio

$$\begin{aligned} \zeta &= \frac{\pi(\tilde{\mathbf{X}}^* | \{\mathbf{y}_i\}_{i=1}^n, \mathbf{U})}{\pi(\tilde{\mathbf{X}}^{(G-1)} | \{\mathbf{y}_i\}_{i=1}^n, \mathbf{U})} \frac{q(\tilde{\mathbf{X}}^{(G-1)} | \tilde{\mathbf{X}}^*, \mathbf{U})}{q(\tilde{\mathbf{X}}^* | \tilde{\mathbf{X}}^{(G-1)}, \mathbf{U})} \\ &= \frac{P(\{\mathbf{y}_i\}_{i=1}^n | \tilde{\mathbf{X}}^*)}{P(\{\mathbf{y}_i\}_{i=1}^n | \tilde{\mathbf{X}}^{(G-1)})} \frac{NB(\{\mathbf{y}_i\}_{i=1}^n | \tilde{\mathbf{X}}^{(G-1)}, \mathbf{R})}{NB(\{\mathbf{y}_i\}_{i=1}^n | \tilde{\mathbf{X}}^*, \mathbf{R})}, \end{aligned}$$

where $\mathbf{R} = \{r_{it}\}$ is matrix for dispersion parameters for each neuron at all time points. $P(\cdot)$ denotes the Poisson likelihood, and $NB(\cdot)$ denotes the negative binomial likelihood. Accept the proposal $\tilde{\mathbf{X}}^*$ with probability $\min(1, \zeta)$.

Table 1: PG-MH vs. Gaussian approximation for μ and X

	μ	X_{*1}	X_{*2}
$\cos(\hat{\xi}_{\text{PGMH}}, \xi_{\text{true}})$	0.9724	0.7663	0.9728
$\cos(\hat{\xi}_{\mathcal{N}}, \xi_{\text{true}})$	0.9680	0.7443	0.9699
$\cos(\hat{\xi}_{\text{PGMH}}, \hat{\xi}_{\mathcal{N}})$	0.9435	0.9664	0.9959

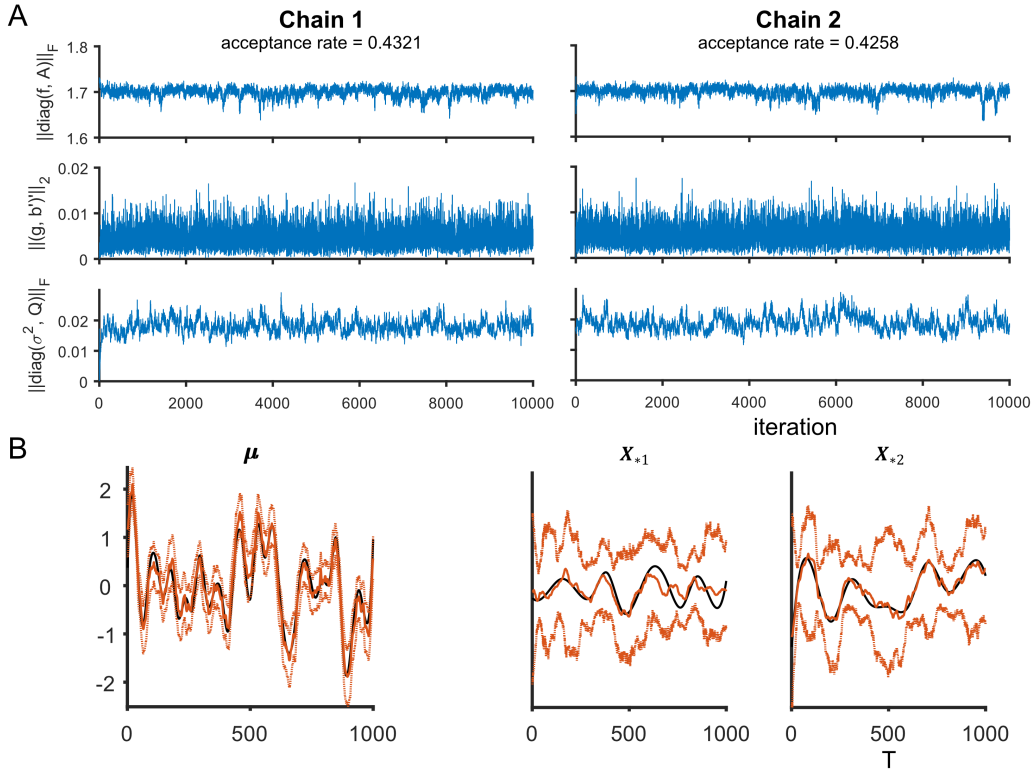


Figure 1: **PG-MH for DPFA** Two independent chains (10,000 iterations) for cluster 1 data in Figure 2. The full conditional distributions of μ and X are sampled by PG-MH algorithm (Algorithm 1). **A.** Traceplot of Frobenius norms of linear dynamics. The acceptance rates in subtitle are for sampling μ and X (PG-MH step). **B.** The true (black) and fitted (colored) population baseline and latent factor. Use samples from iteration 5000 to 10,000 in chain 1, the solid orange lines show averages and the dashed lines show 95% HPD intervals.

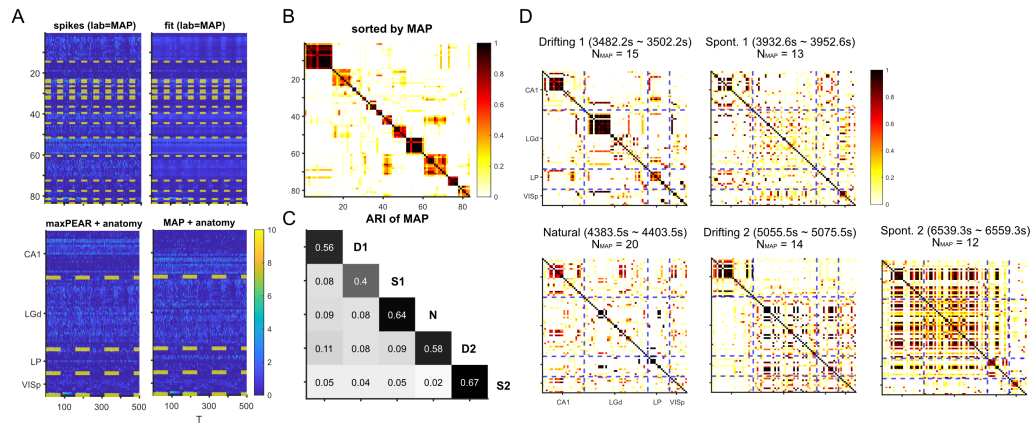


Figure 2: Supplementary Results for Neuropixels Application A. The first row shows the spike counts and fitted mean firing rate, sorted according to the MAP label estimates. The second row sort the spikes further according the anatomical sites, which shows clustering structure within each region. **B.** The posterior similarity matrix sorted by MAP estimates. **C.** The ARI of MAP estimates. The diagonal is ARI between 2 chains, and the off-diagonal is mean ARI of MAP for 4 combinations. **D.** Posterior similarity matrices for the second chains. Neurons are sorted as in the first panel of Fig. 4E

References

- Carter, C. K. and Kohn, R. (1994). On Gibbs sampling for state space models. *Biometrika*, 81(3):541–553.
- Duane, S., Kennedy, A., Pendleton, B. J., and Roweth, D. (1987). Hybrid monte carlo. *Physics Letters B*, 195(2):216–222.
- Eden, U. T., Frank, L. M., Barbieri, R., Solo, V., and Brown, E. N. (2004). Dynamic Analysis of Neural Encoding by Point Process Adaptive Filtering. *Neural Computation*, 16(5):971–998.
- Fokoué, E. and Titterton, D. M. (2003). Mixtures of factor analysers. Bayesian estimation and inference by stochastic simulation. *Machine Learning*, 50(1-2):73–94.
- Frühwirth-Schnatter, S. (1994). Data augmentation and dynamic linear models. *Journal of Time Series Analysis*, 15:183–202.
- Jain, S. and Neal, R. M. (2004). A split-merge markov chain monte carlo procedure for the dirichlet process mixture model. *Journal of Computational and Graphical Statistics*, 13(1):158–182.
- Jain, S. and Neal, R. M. (2007). Splitting and merging components of a nonconjugate Dirichlet process mixture model. *Bayesian Analysis*, 2(3):445–472.
- Krzanowski, W. J. and Marriott, F. H. C. (1994a). Multivariate analysis. Part 2: Classification, covariance structures and repeated measurements.
- Krzanowski, W. J. and Marriott, F. H. C. (1994b). Multivariate analysis. Part I. Distributions, ordination and inference.
- Macke, J. H., Buesing, L., Cunningham, J. P., Yu, B. M., Shenoy, K. V., and Sahani, M. (2011). Empirical models of spiking in neural populations. *Advances in Neural Information Processing Systems*, 24.
- Miller, J. W. and Harrison, M. T. (2018). Mixture models with a prior on the number of components. *Journal of the American Statistical Association*, 113(521):340.
- Neal, R. M. (2000). Markov Chain Sampling Methods for Dirichlet Process Mixture Models. *Journal of Computational and Graphical Statistics*, 9(2):249–265.
- Paninski, L., Ahmadian, Y., Ferreira, D. G., Koyama, S., Rahnama Rad, K., Vidne, M., Vogelstein, J., and Wu, W. (2010). A new look at state-space models for neural data. *Journal of Computational Neuroscience*, 29(1-2):107–126.
- Polson, N. G., Scott, J. G., and Windle, J. (2013). Bayesian inference for logistic models using pólya–gamma latent variables. *Journal of the American Statistical Association*, 108(504):1339–1349.
- Rauch, H. E., Tung, F., and Striebel, C. T. (1965). Maximum likelihood estimates of linear dynamic systems. *AIAA Journal*, 3(8):1445–1450.
- Sen, D., Patra, S., and Dunson, D. (2018). Constrained inference through posterior projections.
- Windle, J., Carvalho, C. M., Scott, J. G., and Sun, L. (2013). Efficient Data Augmentation in Dynamic Models for Binary and Count Data.