

# Model-based Clustering for Neural Populations

The goal for this research is to do model-based clustering for neural populations, by making use of features for each counting process observation.

## 1 Notations

Assume we can observe neural activities for  $N$  neurons, with counting observation up to  $T$  steps. Therefore, the observation is a  $N$ -by- $T$  matrix,  $\mathbf{Y} \in \mathbb{Z}_{\geq 0}^{N \times T}$ , with each row represents the recording from single neuron. Denote the recording for neuron  $i$  as  $\mathbf{y}_i = (y_{i1}, \dots, y_{iT})'$ ,  $i = 1, \dots, N$ , with the cluster index for neuron  $i$  as  $z_i \in \{1, \dots\}$ . The number of neurons in cluster  $j$  is  $n_j = \sum_{i=1}^N I(z_i = j)$ , and  $\sum_{j=1,2,\dots} n_j = N$ . The proportion/ probability in cluster  $z_i$  is  $\rho_{z_i}$ .

## 2 Clustering Wrapper

The model-based clustering problem can be transformed into fitting the mixture model (MM). The likelihood for each cluster depends on how we model the counting observation, but fitting strategies for MM are the same for all models. Here, I choose to fit the MM by Gibbs sampler. Depending on whether the number of cluster is finite or not, there are two versions: finite mixture model (FMM) and Dirichlet process mixture model (DPMM).

### 2.1 Finite Mixture Model

Assume the number of cluster is  $J$ . The full likelihood for these  $N$  neurons is

$$L = \prod_{i=1}^N \rho_{z_i} f(\mathbf{y}_i | \boldsymbol{\Theta}_{z_i}) = \prod_{j=1}^J \rho_j^{n_j} \left[ \prod_{i: z_i=j} f(\mathbf{y}_i | \boldsymbol{\Theta}_j) \right]$$

, where  $\boldsymbol{\Theta}_j$  contains all parameters in cluster  $j$  defined by the specific model. Therefore, the parameters need to update are:

- (1) Cluster indicator:  $\{z_i\}_{i=1}^N$
- (2) Cluster proportion:  $\rho = (\rho_1, \dots, \rho_J)'$
- (3) Model parameters:  $\boldsymbol{\Theta}_j$

The (conditional) priors for clustering-related parameters:

- (1) Cluster indicator  $\{z_i\}_{i=1}^N$ :  $P(z_i = j) = \rho_j$
- (2) Cluster proportion  $\rho = (\rho_1, \dots, \rho_J)'$ :

$$\rho \sim \text{Dir}(\delta_1, \dots, \delta_J)$$

, where  $\delta_1 = \dots = \delta_J = 1$

So, the MCMC( Gibbs sampler) iteration for FMM is:

- (1) Update  $\{z_i\}_{i=1}^N$ :

$$P(z_i = j | \mathbf{y}_i, \{\boldsymbol{\Theta}_j\}_{j=1}^J) \propto \rho_j f(\mathbf{y}_i | \boldsymbol{\Theta}_j)$$

- (2) Update  $\rho = (\rho_1, \dots, \rho_J)'$ :

$$\rho | \{\mathbf{y}_i\}_{i=1}^N, \{z_i\}_{i=1}^N, \{\boldsymbol{\Theta}_j\}_{j=1}^J \sim \text{Dir}(\delta_1 + n_1, \dots, \delta_J + n_J)$$

- (3) Update  $\boldsymbol{\Theta}_j$ : this is defined by the specific model. When there's no  $z_i = j$ , just sample  $\boldsymbol{\Theta}_j$  from priors or by other observation-independent ways.

## 2.2 Dirichlet Process Mixture Model

Since calculation of posterior predictive distribution can be hard or even impossible for complicated models, instead of using the popular CRP representation of DP (Neal, 2020), I choose to use the slice sampler (Walker, 2007).

Use the "stick-breaking" construction for cluster proportion, i.e.

$$\rho_1 = \eta_1$$

$$\rho_j = (1 - \eta_1) \cdot \dots \cdot (1 - \eta_{j-1}) \eta_j$$

$$\eta_j \sim \text{Beta}(1, \alpha)$$

In the slice sampler for DPMM, the parameters need to update are:

- (1) "stick-breaking" elements:  $\eta_j$
- (2) Augment latent variable:  $\{u_i\}_{i=1}^N$
- (3) Model parameters:  $\boldsymbol{\Theta}_j$
- (4) Cluster indicator:  $\{z_i\}_{i=1}^N$

So, the MCMC( Gibbs sampler) iteration for DPMM is:

(1) update  $\eta_j$ , for  $j = 1, \dots, z^* = \max \{z_i\}_{i=1}^N$  as

$$\eta_j | \{z_i\}_{i=1}^N, \dots \sim \text{Beta}(n_j + 1, N - \sum_{l=1}^j n_l + \alpha)$$

(2) update  $\{u_i\}_{i=1}^N$ :

$$u_i | \rho, \dots \sim U(0, \rho_{z_i})$$

(3) update  $\eta_j$ , for  $j = z^* + 1, \dots, s^*$ .  $s^*$  is the smallest value, s.t.  $\sum_{j=1}^{s^*} \rho_j > 1 - \min \{u_1, \dots, u_N\}$

$$\eta_j \sim \text{Beta}(1, \alpha)$$

(4) Update  $\Theta_j$ : this is defined by the specific model. When there's no  $z_i = j$ , just sample  $\Theta_j$  from priors or by other observation-independent ways.

(5) Update  $\{z_i\}_{i=1}^N$

$$P\left(z_i = j \middle| \mathbf{y}_i, \{\Theta_j\}, \rho, \{u_i\}_{i=1}^N\right) = \frac{f(\mathbf{y}_i | \Theta_j)}{\sum_{j: \rho_j > u_i} f(\mathbf{y}_i | \Theta_j)}$$

### 3 linear Dynamical System Model

Here, I model the observations by a linear dynamical system (LDS) model.

LDS models the multi-dimensional time series using a lower dimensional latent representation of the system, which evolves over time according to linear dynamics. By specifying the linear dynamics and process noise covariance, we can also handle the interactions between different neural populations (clusters).

#### 3.1 Model Details

Denote the latent vector in cluster  $j$  as  $\mathbf{x}_t^{(j)} \in \mathbb{R}^{p_j}$ . For simplicity, assume all  $p_j = p$ . Each observation follows a Poisson distribution:

$$\log \lambda_{it} = d_i + \mathbf{c}'_i \mathbf{x}_t^{(z_i)}$$

$$y_{it} \sim \text{Poisson}(\lambda_{it})$$

, where  $\mathbf{c}_i \in \mathbb{R}^p$  and  $\mathbf{x}_t^{(z_i)} \in \mathbb{R}^p$ .

Although the loading ( $d_i$  and  $\mathbf{c}_i$ ) is determined by neuron index  $i$ , the distribution is also cluster-dependent. That is,

$$(d_i, \mathbf{c}'_i)' \sim N(\boldsymbol{\mu}_{\text{dc}}^{(z_i)}, \boldsymbol{\Sigma}_{\text{dc}}^{(z_i)})$$

By doing this, the loading within each cluster is also correlated.

Denote all latent states as  $\mathbf{x}_t = \left( \mathbf{x}_t^{(1)}, \mathbf{x}_t^{(2)}, \dots \right)'$  and they evolve linearly with a Gaussian noise:

$$\begin{aligned}\mathbf{x}_1 &\sim N(\mathbf{x}_0, \mathbf{Q}_0) \\ \mathbf{x}_{t+1}|\mathbf{x}_t &\sim N(\mathbf{A}\mathbf{x}_t + \mathbf{b}, \mathbf{Q})\end{aligned}$$

For simplicity, assume  $\mathbf{Q}_0$  is known (e.g.  $\mathbf{Q}_0 = \mathbf{I} \times 10^{-2}$ ).

If we assume process noise covariance is block diagonal (Joshua et al., 2020), we can write things as:

$$\mathbf{x}_{t+1}^{(j)}|\mathbf{x}_t^{(1)}, \mathbf{x}_t^{(2)}, \dots \sim N\left(\sum_{l=1, \dots} \mathbf{A}_{j \leftarrow l} \mathbf{x}_t^{(l)} + \mathbf{b}_j, \mathbf{Q}^{(j)}\right)$$

Notice  $\{\mathbf{A}_{j \leftarrow l}\}$  forms the full transition matrix as:

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{1 \leftarrow 1} & \mathbf{A}_{1 \leftarrow 2} & \dots \\ \mathbf{A}_{2 \leftarrow 1} & \mathbf{A}_{2 \leftarrow 2} & \dots \\ \dots & \dots & \dots \end{pmatrix}$$

Denote the  $j^{\text{th}}$  row block of  $\mathbf{A}$  as  $\mathbf{A}_j = (\mathbf{A}_{j \leftarrow 1} \quad \mathbf{A}_{j \leftarrow 2} \quad \dots)$ . Then,  $\sum_{l=1, \dots} \mathbf{A}_{j \leftarrow l} \mathbf{x}_t^{(l)} + \mathbf{b}_j = \mathbf{A}_j \mathbf{x}_t + \mathbf{b}_j$ .

If we further let  $\mathbf{Q}$  be diagonal, with the  $k^{\text{th}}$  row of  $\mathbf{x}_t$ ,  $\mathbf{A}$ ,  $\mathbf{b}$  denoted as  $x_{kt}$ ,  $\mathbf{a}_k$ ,  $b_k$ . The corresponding process noise variance is  $q_k$ . Then:

$$x_{k,t+1}|x_{kt} \sim N\left(\mathbf{a}_k' \mathbf{x}_t + b_k, q_k\right)$$

To facilitate derivation in Gibbs sampler, write the linear dynamics of  $\mathbf{x}_t$  in MV-GLM form, i.e.

$$\mathbf{x}_{t+1}' = \begin{pmatrix} 1 & \mathbf{x}_t' \end{pmatrix} \begin{pmatrix} \mathbf{b}' \\ \mathbf{A}' \end{pmatrix} + \epsilon_t'$$

, where  $\epsilon \sim N(0, \mathbf{Q})$ . Then if we stack the latent by row, the problem is reduced to a multivariate general linear model (MV-GLM) problem:

$$\begin{pmatrix} \mathbf{x}_2' \\ \mathbf{x}_3' \\ \vdots \\ \mathbf{x}_T' \end{pmatrix} = \begin{pmatrix} 1 & \mathbf{x}_1' \\ 1 & \mathbf{x}_2' \\ \vdots & \vdots \\ 1 & \mathbf{x}_{T-1}' \end{pmatrix} \begin{pmatrix} \mathbf{b}' \\ \mathbf{A}' \end{pmatrix} + \mathbf{E}$$

, where  $\mathbf{E} = (\epsilon_1, \dots, \epsilon_{T-1})'$ .

In summary, Let  $\Theta = \left\{ \mathbf{z}, \mathbf{d}, \mathbf{C}, \{\mathbf{x}_t\}_{t=1}^T, \mathbf{x}_0, \mathbf{Q}_0, \mathbf{A}, \mathbf{b}, \mathbf{Q} \right\}$  be the set of parameters for LDS. The number of observation/ neuron is  $N$ , and they can group

into  $J$  clusters. In each cluster, there are  $p$  latent state vectors. The recording length is  $T$ .

- (1)  $\mathbf{z} \in \mathbb{Z}^N$ : cluster index for each neuron/ observation, with  $z_i \in \{1, \dots, J\}$  for  $i = 1, \dots, N$ .
- (2)  $\mathbf{d} \in \mathbb{R}^N$  and  $\mathbf{C} \in \mathbb{R}^{N \times p}$ : baseline & loading of the latent
- (3)  $\mathbf{x}_t \in \mathbb{R}^{Jp}$ : all latent at  $t$ .
- (4)  $\mathbf{x}_t^{(j)} \in \mathbb{R}^p$  means the latent in cluster  $j$  at time  $t$ .
- (5)  $\mathbf{x}_0 \in \mathbb{R}^{Jp}$  and  $\mathbf{Q}_0 \in \mathbb{R}^{Jp \times Jp}$ : mean and covariance of  $\mathbf{x}_1$ . For simplicity, assume  $\mathbf{Q}_0$  is known.
- (6)  $\mathbf{A} \in \mathbb{R}^{Jp \times Jp}$ ,  $\mathbf{b} \in \mathbb{R}^{Jp}$ ,  $\mathbf{Q} \in \mathbb{R}^{Jp \times Jp}$ : linear dynamics of the latent.

Since the progress noise is independent at each step and the observation is assumed conditional independent, the likelihood is:

$$f(\mathbf{y}|\Theta) = \prod_{i=1}^N \prod_{t=1}^T P(y_{it}|\Theta) = \prod_{i=1}^N \prod_{t=1}^T POI(y_{it}|\exp(d_i + \mathbf{c}'_i \mathbf{x}_t^{(z_i)}))$$

, where  $POI(\cdot|\lambda)$  is the density of  $Poisson(\lambda)$ .

#### DETOUR:

Maybe in the future, we may switch to EM. To help with this, I also give the likelihood for complete observation  $(\mathbf{y}, \{\mathbf{x}_t\}_{t=1}^T)$ . Let  $\Theta' = \Theta \setminus \{\mathbf{x}_t\}_{t=1}^T$ .

$$f(\mathbf{y}, \{\mathbf{x}_t\}_{t=1}^T | \Theta') = \prod_{i=1}^N \prod_{t=1}^T P(y_{it}, \mathbf{x}_t | \Theta') = \prod_{i=1}^N \prod_{t=1}^T POI(y_{it}|\exp(d_i + \mathbf{c}'_i \mathbf{x}_t^{(z_i)})) \cdot N(\mathbf{x}_t^{(z_i)} | \mathbf{A}_{z_i} \mathbf{x}_t + \mathbf{b}_{z_i}, \mathbf{Q}_{z_i})$$

, where  $\mathbf{A}_j \in \mathbb{R}^{p \times Jp}$  and  $\mathbf{b}_j \in \mathbb{R}^p$  are the  $j^{\text{th}}$  row block of  $\mathbf{A}$  and  $\mathbf{b}$ .  $\mathbf{Q}_j \in \mathbb{R}^{p \times p}$  is the  $j^{\text{th}}$  row and column block of  $\mathbf{Q}$ .

### 3.2 Conditional Priors for Parameters

The parameters need to estimate:

- (1) Latent vectors:  $\{\mathbf{x}_t\}_{t=1}^T$
- (2) Initials:  $\mathbf{x}_0$
- (3) Linear mapping (loading) for latent vectors:  $\{d_i\}_{i=1}^N$  and  $\{\mathbf{c}_i\}_{i=1}^N$
- (4) Mean and covariance for loading in each cluster:  $\{\boldsymbol{\mu}_{\text{dc}}^{(j)}\}_j$  and  $\{\boldsymbol{\Sigma}_{\text{dc}}^{(j)}\}_j$
- (5) Linear dynamics for latent vectors:  $\mathbf{A}$  and  $\mathbf{b}$
- (6) Process noise:  $\mathbf{Q}$

The conditional priors for these parameters:

- (1) Latent vectors  $\{\mathbf{x}_t\}_{t=1}^T$ : the conditional prior is defined by

$$\mathbf{x}_1 \sim N(\mathbf{x}_0, \mathbf{Q}_0)$$

$$\mathbf{x}_{t+1}|\mathbf{x}_t \sim N(\mathbf{A}\mathbf{x}_t + \mathbf{b}, \mathbf{Q})$$

- (2) Initials  $\mathbf{x}_0$ : assume there are  $J$  clusters,

$$\mathbf{x}_0 \sim N(\boldsymbol{\mu}_{\mathbf{x}_{00}}, \boldsymbol{\Sigma}_{\mathbf{x}_{00}})$$

, where  $\boldsymbol{\mu}_{\mathbf{x}_{00}} = \mathbf{0}_{Jp}$  and  $\boldsymbol{\Sigma}_{\mathbf{x}_{00}} = \mathbf{I}_{Jp}$

- (3) Linear mapping (loading) for latent vectors  $\{d_i\}_{i=1}^N$  and  $\{\mathbf{c}_i\}_{i=1}^N$ :

$$(d_i, \mathbf{c}_i')' \sim N(\boldsymbol{\mu}_{\text{dc}}^{(z_i)}, \boldsymbol{\Sigma}_{\text{dc}}^{(z_i)})$$

- (4) Mean and covariance for loading in each cluster Mean and covariance for loading in each cluster  $\{\boldsymbol{\mu}_{\text{dc}}^{(j)}\}_j$  and  $\{\boldsymbol{\Sigma}_{\text{dc}}^{(j)}\}_j$ :

$$\boldsymbol{\mu}_{\text{dc}}^{(j)} \sim N(\delta_{dc0}, \mathbf{T}_{dc0})$$

, where  $\delta_{dc0} = \mathbf{0}_{p+1}$  and  $\mathbf{T}_{dc0} = \mathbf{I}_{p+1}$

$$\boldsymbol{\Sigma}_{\text{dc}}^{(j)} \sim W^{-1}(\Psi_{dc0}, \nu_{dc0})$$

, where  $\nu_{dc0} = p + 1 + 2$  and  $\Psi_{dc0} = \mathbf{I}_{p+1} \times 10^{-4}$

- (5) Linear dynamics for latent vectors:  $\mathbf{A}$ ,  $\mathbf{b}$  and process noise  $\mathbf{Q}$ .

Denote  $\mathbf{F} = \begin{pmatrix} \mathbf{b}' \\ \mathbf{A} \end{pmatrix}$ ,  $\mathbf{f} = \text{vec}(\mathbf{F})$

$$P(\mathbf{f}, \mathbf{Q}) = P(\mathbf{Q})P(\mathbf{f}|\mathbf{Q})$$

$$\mathbf{Q} \sim W^{-1}(\Psi_{\mathbf{Q}_0}, \nu_{\mathbf{Q}_0})$$

$$\mathbf{f} \sim N(\mathbf{f}_0, \mathbf{Q} \otimes \boldsymbol{\Lambda}_0^{-1})$$

, where  $\mathbf{f}_0 = \text{vec}(\mathbf{F}_0)$ . If the number of cluster is  $J$ ,  $\nu_{\mathbf{Q}_0} = Jp + 2$ ,  $\Psi_{\mathbf{Q}_0} = \mathbf{I}_{Jp} \times 10^{-4}$ . (To make the mean of  $\mathbf{Q}$  loosely centered around  $\mathbf{I}_{Jp} \times 10^{-4}$ ),

$\mathbf{F}_0 = \begin{pmatrix} \mathbf{0}_{Jp}' \\ \mathbf{I}_{Jp} \end{pmatrix}$  and  $\boldsymbol{\Lambda}_0 = \mathbf{I}_{Jp+1}$

### 3.3 MCMC (Gibbs Sampler)

#### 3.3.1 Update $\{\mathbf{x}_t\}_{t=1}^T$

Use Laplace approximation and make use of the block tri-diagonal Hessian.

Denote  $t^{\text{th}}$  column of mean firing rate and observation as  $\tilde{\boldsymbol{\lambda}}_t = (\lambda_{1t}, \dots, \lambda_{Nt})'$  and  $\tilde{\mathbf{y}}_t = (y_{1t}, \dots, y_{Nt})'$ . The linear mapping matrix for all observations is  $\mathbf{C}$ , such that  $\log \tilde{\boldsymbol{\lambda}}_t = \mathbf{d} + \mathbf{C}\mathbf{x}_t$ . Let  $\mathbf{x} = (\mathbf{x}'_1, \dots, \mathbf{x}'_T)'$  and  $f(\mathbf{x}) = \log P(\mathbf{x} | \{\mathbf{y}_i\}_{i=1}^N, \mathbf{C}, \mathbf{Q}_0, \mathbf{A}, \mathbf{b}, \mathbf{Q}, \dots)$ . The first and second derivative with respect to  $\mathbf{x}$ , for  $t = 2, \dots, T-1$ :

$$\begin{aligned} \frac{\partial f}{\partial \mathbf{x}_1} &= \mathbf{C}' (\tilde{\mathbf{y}}_1 - \tilde{\boldsymbol{\lambda}}_1) - \mathbf{Q}_0^{-1} (\mathbf{x}_1 - \mathbf{x}_0) + \mathbf{A}' \mathbf{Q}^{-1} (\mathbf{x}_2 - \mathbf{A}\mathbf{x}_1 - \mathbf{b}) \\ \frac{\partial f}{\partial \mathbf{x}_t} &= \mathbf{C}' (\tilde{\mathbf{y}}_t - \tilde{\boldsymbol{\lambda}}_t) - \mathbf{Q}^{-1} (\mathbf{x}_t - \mathbf{A}\mathbf{x}_{t-1} - \mathbf{b}) + \mathbf{A}' \mathbf{Q}^{-1} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t - \mathbf{b}) \\ \frac{\partial f}{\partial \mathbf{x}_T} &= \mathbf{C}' (\tilde{\mathbf{y}}_T - \tilde{\boldsymbol{\lambda}}_T) - \mathbf{Q}^{-1} (\mathbf{x}_T - \mathbf{A}\mathbf{x}_{T-1} - \mathbf{b}) \\ \frac{\partial^2 f}{\partial \mathbf{x}_1 \partial \mathbf{x}'_1} &= -\mathbf{C}' \text{Diag}(\tilde{\boldsymbol{\lambda}}_1) \mathbf{C} - \mathbf{Q}_0^{-1} - \mathbf{A}' \mathbf{Q}^{-1} \mathbf{A} \\ \frac{\partial^2 f}{\partial \mathbf{x}_t \partial \mathbf{x}'_t} &= -\mathbf{C}' \text{Diag}(\tilde{\boldsymbol{\lambda}}_t) \mathbf{C} - \mathbf{Q}^{-1} - \mathbf{A}' \mathbf{Q}^{-1} \mathbf{A} \\ \frac{\partial^2 f}{\partial \mathbf{x}_T \partial \mathbf{x}'_T} &= -\mathbf{C}' \text{Diag}(\tilde{\boldsymbol{\lambda}}_T) \mathbf{C} - \mathbf{Q}^{-1} \\ \frac{\partial^2 f}{\partial \mathbf{x}_1 \partial \mathbf{x}'_2} &= \frac{\partial^2 f}{\partial \mathbf{x}_t \partial \mathbf{x}'_{t+1}} = \mathbf{A}' \mathbf{Q}^{-1} \end{aligned} \quad \frac{\partial^2 f}{\partial \mathbf{x}_t \partial \mathbf{x}'_{t-1}} = \mathbf{Q}^{-1} \mathbf{A}$$

So, the gradient is:

$$\nabla = \frac{\partial f}{\partial \mathbf{x}} = \left( \left( \frac{\partial f}{\partial \mathbf{x}_1} \right)', \dots, \left( \frac{\partial f}{\partial \mathbf{x}_T} \right)' \right)'$$

And the block tri-diagonal Hessian:

$$H = \frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{x}'} = \begin{pmatrix} \frac{\partial^2 f}{\partial \mathbf{x}_1 \partial \mathbf{x}'_1} & \mathbf{A}' \mathbf{Q}^{-1} & 0 & \dots & 0 \\ \mathbf{Q}^{-1} \mathbf{A} & \frac{\partial^2 f}{\partial \mathbf{x}_2 \partial \mathbf{x}'_2} & \mathbf{A}' \mathbf{Q}^{-1} & \dots & \vdots \\ 0 & \mathbf{Q}^{-1} \mathbf{A} & \frac{\partial^2 f}{\partial \mathbf{x}_3 \partial \mathbf{x}'_3} & \dots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & \dots & \frac{\partial^2 f}{\partial \mathbf{x}_T \partial \mathbf{x}'_T} \end{pmatrix}$$

Use Newton-Raphson to find  $\boldsymbol{\mu}_{\mathbf{x}} = \text{argmax}_{\mathbf{x}} (f(\mathbf{x}))$  and  $\boldsymbol{\Sigma}_{\mathbf{x}} = - \left[ \frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{x}'} \Big|_{\mathbf{x}=\boldsymbol{\mu}_{\mathbf{x}}} \right]^{-1}$ , such that  $(P(\mathbf{x} | \{\mathbf{y}_i\}_{i=1}^N, \dots) \approx N(\boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{x}})$ . When using Newton-Raphson

(NR),  $H \setminus \nabla$  in MATLAB will make use of block tri-diagonal structure automatically.

However, NR is not robust to bad initials. At the first few iterations, simply using fitting from previous step may lead to infinite Hessian. When the initial from previous step fails, use the approximation at recursive priors, i.e. the adaptive smoother estimates, as the initial. The adaptive smoother estimates are from backward RTS smoother from adaptive filter, and the details about Poisson adaptive filter can be found in Eden et al., 2004.

To sample efficiently and make best use of sparse covariance, use Cholesky decomposition of  $\Sigma_{\mathbf{x}}^{-1} = \mathbf{R}\mathbf{R}'$ : sample  $\mathbf{Z} \sim N(\mathbf{R}'\mu_{\mathbf{x}}, \mathbf{I})$ , then  $\mathbf{x} = (\mathbf{R}')^{-1}\mathbf{Z} \sim N(\mu_{\mathbf{x}}, \Sigma_{\mathbf{x}})$ .

#### Step Further:

Using normal approximation may not be accurate enough, and this may lead to a bad mixing. We can step further to use the Metropolis-Hasting, based on the normal approximated variance. To be specific, the proposal distribution is  $Q(\mathbf{x}^*|\mathbf{x}^{(s)}) \sim N(\mathbf{x}^{(s)}, \alpha\Sigma_{\mathbf{x}})$ , where  $\alpha$  is a scalar to make the acceptance rate to be 0.4 to 0.6. According to Roberts and Rosenthal, 2001. For high dimensional MH, the optimal proposal is  $N(x, 2.38^2\Sigma/d)$ . See details & experiments in Section 4: Simulations.

#### 3.3.2 Update $\mathbf{x}_0$

$$P(\mathbf{x}_0|\mathbf{x}_1, \mathbf{Q}_0 \dots) \propto N(\mathbf{x}_1|\mathbf{x}_0, \mathbf{Q}_0)N(\mathbf{x}_0|\mu_{\mathbf{x}_0}, \Sigma_{\mathbf{x}_0})$$

By conjugacy,  $\mathbf{x}_0|\mathbf{x}_1, \mathbf{Q}_0 \dots \sim N(\mu_{\mathbf{x}_0}, \Sigma_{\mathbf{x}_0})$

$$\begin{aligned}\Sigma_{\mathbf{x}_0} &= [\Sigma_{\mathbf{x}_{00}}^{-1} + \mathbf{Q}_0^{-1}]^{-1} \\ \mu_{\mathbf{x}_0} &= \Sigma_{\mathbf{x}_0} (\Sigma_{\mathbf{x}_{00}}^{-1}\mu_{\mathbf{x}_{00}} + \mathbf{Q}_0^{-1}\mathbf{x}_1)\end{aligned}$$

#### 3.3.3 Update $\{d_i\}_{i=1}^N$ and $\{\mathbf{c}_i\}_{i=1}^N$

To update efficiently, use Laplace approximation and the error is independent for each row (i.e. update things row by row). Denote  $(d_i, \mathbf{c}_i')' = \zeta_i \in \mathbb{R}^{p+1}$  and  $(1, \mathbf{x}_t'^{(z_i)}) = \tilde{\mathbf{x}}_t'^{(z_i)}$ .

$$P\left(\zeta_i \middle| \mathbf{y}_i, \left\{\mathbf{x}_t^{(z_i)}\right\}_{t=1}^T, \dots\right) = \exp f(\zeta_i) \approx N(\zeta_i | \mu_{\zeta_i}, \Sigma_{\zeta_i})$$

$$\begin{aligned}\frac{\partial f}{\partial \zeta_i} &= \frac{\partial l}{\partial \zeta_i} - \Sigma_{\text{dc}}^{(z_i)-1} \left( \zeta_i - \mu_{\text{dc}}^{(z_i)} \right) = \left[ \sum_{t=1}^T \tilde{\mathbf{x}}_t^{(z_i)} (y_{it} - \lambda_{it}) \right] - \Sigma_{\text{dc}}^{(z_i)-1} \left( \zeta_i - \mu_{\text{dc}}^{(z_i)} \right) \\ \frac{\partial^2 f}{\partial \zeta_i \partial \zeta_i'} &= \frac{\partial^2 l}{\partial \zeta_i \partial \zeta_i'} - \Sigma_{\text{dc}}^{(z_i)-1} = - \left[ \sum_{t=1}^T \lambda_{it} \tilde{\mathbf{x}}_t^{(z_i)} \tilde{\mathbf{x}}_t'^{(z_i)} \right] - \Sigma_{\text{dc}}^{(z_i)-1}\end{aligned}$$



, where  $l$  is Poisson log-likelihood. Then use Newton-Raphson to find  $\boldsymbol{\mu}_{\zeta_i} = \text{argmax}_{\zeta_i} (f(\zeta_i))$  and  $\boldsymbol{\Sigma}_{\zeta_i} = - \left[ \frac{\partial^2 f}{\partial \zeta_i \partial \zeta_i'} \Big|_{\zeta_i = \boldsymbol{\mu}_{\zeta_i}} \right]^{-1}$ . If initial as the previous step fits fails, simply use prior mean of  $\boldsymbol{\mu}_{\zeta_i}$ , i.e.  $\boldsymbol{\delta}_{dc0}$ .

**Step further:** Again, use MH to sample the posterior, based on the normal approximated variance. Since the dimension of  $d$  and  $C$  is not large (just  $d = 3$ ) in the simulation, I didn't use the optimal scalar yet. **Modify it later.**

### 3.3.4 Update $\{\boldsymbol{\mu}_{dc}^{(j)}\}_j$ and $\{\boldsymbol{\Sigma}_{dc}^{(j)}\}_j$

Purpose: to make loading within each cluster depends on each other, and this will help with clustering. As above, denote  $(d_i, \mathbf{c}_i')' = \boldsymbol{\zeta}_i \in \mathbb{R}^{p+1}$ .

- (1) Mean  $\{\boldsymbol{\mu}_{dc}^{(j)}\}_j$ : by conjugacy,  $\boldsymbol{\mu}_{dc}^{(j)} \sim N(\boldsymbol{\delta}_{dc}, \mathbf{T}_{dc})$

$$\mathbf{T}_{dc}^{-1} = \left( \mathbf{T}_{dc0}^{-1} + n_j \boldsymbol{\Sigma}_{dc}^{(j)-1} \right)^{-1}$$

$$\boldsymbol{\delta}_{dc} = \mathbf{T}_{dc} \left( \mathbf{T}_{dc0}^{-1} \boldsymbol{\delta}_{dc0} + \boldsymbol{\Sigma}_{dc}^{(j)-1} \sum_{i:z_i=j} \boldsymbol{\zeta}_i \right)$$

- (2) Covariance  $\{\boldsymbol{\Sigma}_{dc}^{(j)}\}_j$ : by conjugacy,  $\boldsymbol{\Sigma}_{dc}^{(j)} \sim W^{-1}(\Psi_{dc}, \nu_{dc})$

$$\nu_{dc} = n_j + \nu_{dc0}$$

$$\Psi_{dc} = \Psi_{dc0} + \sum_{i:z_i=j} \left( \boldsymbol{\zeta}_i - \boldsymbol{\mu}_{dc}^{(j)} \right) \left( \boldsymbol{\zeta}_i - \boldsymbol{\mu}_{dc}^{(j)} \right)'$$

### 3.3.5 Update $\mathbf{A}$ , $\mathbf{b}$ and $\mathbf{Q}$

Here, I only give the update for full  $\mathbf{Q}$ . If we want to assume block-diagonal or diagonal structure of  $\mathbf{Q}$ , just update things cluster-by-cluster or latent-by-latent.

As shown before, the problem is the usual Bayesian MV-GLM problem. Denote

$$\mathbf{F} = \begin{pmatrix} \mathbf{b}' \\ \mathbf{A}' \end{pmatrix}, \mathbf{f} = \text{vec}(\mathbf{F}), \mathbf{Y}_{bA} = (\mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_T)' \text{ and } \mathbf{X}_{bA} = \begin{pmatrix} 1 & \mathbf{x}_1' \\ 1 & \mathbf{x}_2' \\ \vdots & \vdots \\ 1 & \mathbf{x}_{T-1}' \end{pmatrix}$$

$$\mathbf{Y}_{bA} = \mathbf{X}_{bA} \mathbf{F} + \mathbf{E}$$

, where  $\mathbf{E} = (\epsilon_1, \dots, \epsilon_{T-1})'$  and  $\epsilon \sim N(0, \mathbf{Q})$ . Then posteriors are

$$\mathbf{Q} | \mathbf{Y}_{bA}, \mathbf{X}_{bA}, \dots \sim W^{-1}(\Psi_{\mathbf{Q}}, \nu_{\mathbf{Q}})$$

$$\mathbf{f} | \mathbf{Y}_{bA}, \mathbf{X}_{bA}, \mathbf{Q}, \dots \sim N(\text{vec}(\mathbf{F}_{bA}), \mathbf{Q} \otimes \Lambda_{bA}^{-1})$$

, with parameters be:

$$\begin{aligned}\Psi_{\mathbf{Q}} &= \Psi_{\mathbf{Q}_0} + (\mathbf{Y}_{\mathbf{bA}} - \mathbf{X}_{\mathbf{bA}}\mathbf{F}_{\mathbf{bA}})'(\mathbf{Y}_{\mathbf{bA}} - \mathbf{X}_{\mathbf{bA}}\mathbf{F}_{\mathbf{bA}}) + (\mathbf{F}_{\mathbf{bA}} - \mathbf{F}_0)'\mathbf{\Lambda}_0(\mathbf{F}_{\mathbf{bA}} - \mathbf{F}_0) \\ \nu_{\mathbf{Q}} &= \nu_{\mathbf{Q}_0} + T - 1 \\ \mathbf{F}_{\mathbf{bA}} &= (\mathbf{X}_{\mathbf{bA}}'\mathbf{X}_{\mathbf{bA}} + \mathbf{\Lambda}_0)^{-1}(\mathbf{X}_{\mathbf{bA}}'\mathbf{Y}_{\mathbf{bA}} + \mathbf{\Lambda}_0\mathbf{F}_0) \\ \mathbf{\Lambda}_{\mathbf{bA}} &= \mathbf{X}_{\mathbf{bA}}'\mathbf{X}_{\mathbf{bA}} + \mathbf{\Lambda}_0\end{aligned}$$

## 4 Simulations

**All the simulation results are old and problematic. Update them later.**

## 5 Problem from Affine Transformation

Let me review the model first.

Assume there are  $J$  clusters, with  $n$  neurons in each cluster. Each cluster is governed by  $p$  latent vectors.  $\boldsymbol{\lambda}_t^{(j)} \in \mathbb{R}^n$ ,  $\mathbf{x}_t^{(j)} \in \mathbb{R}^p$ ,  $\mathbf{d}^{(j)} \in \mathbb{R}^n$  and  $\mathbf{C}^{(j)} \in \mathbb{R}^{n \times p}$ , with  $j = 1, \dots, J$ , be the mean firing rate, latent state, intercept and loading (at time  $t$ ) for each cluster. Then I just stack all the thing together, and denote corresponding things as  $\boldsymbol{\lambda}_t \in \mathbb{R}^{Jn}$ ,  $\mathbf{x}_t \in \mathbb{R}^{Jp}$ ,  $\mathbf{d} \in \mathbb{R}^{Jn}$  and  $\mathbf{C} \in \mathbb{R}^{Jn \times Jp}$ . Notice  $\mathbf{C}$  is block diagonal.

$$\log \begin{pmatrix} \boldsymbol{\lambda}_t^{(1)} \\ \vdots \\ \boldsymbol{\lambda}_t^{(J)} \end{pmatrix} = \log \boldsymbol{\lambda}_t = \begin{pmatrix} \mathbf{d}^{(1)} \\ \vdots \\ \mathbf{d}^{(J)} \end{pmatrix} + \begin{pmatrix} \mathbf{C}^{(1)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \mathbf{C}^{(J)} \end{pmatrix} \begin{pmatrix} \mathbf{x}_t^{(1)} \\ \vdots \\ \mathbf{x}_t^{(J)} \end{pmatrix} = \mathbf{d} + \mathbf{C}\mathbf{x}_t$$

The latent states progress linearly with a Gaussian noise:

$$\mathbf{x}_{t+1}|\mathbf{x}_t \sim N(\mathbf{A}\mathbf{x}_t + \mathbf{b}, \mathbf{Q})$$

, where  $\mathbf{A} \in \mathbb{R}^{Jp \times Jp}$ ,  $\mathbf{b} \in \mathbb{R}^{Jp}$  and  $\mathbf{Q} \in \mathbb{R}^{Jp \times Jp}$ .

Now assume we take an affine transformation of latent states as  $\mathbf{x}_t^* = \mathbf{M}\mathbf{x}_t + \mathbf{g}$ . To preserve the block diagonal structure of  $\mathbf{C}$ ,  $\mathbf{M}$  is also block diagonal. Then,

$$\begin{aligned}\log \boldsymbol{\lambda}_t &= \mathbf{d} + \mathbf{C}\mathbf{M}^{-1}(\mathbf{M}\mathbf{x}_t + \mathbf{g}) - \mathbf{C}\mathbf{M}^{-1}\mathbf{g} \\ &= (\mathbf{d} - \mathbf{C}\mathbf{M}^{-1}\mathbf{g}) + \mathbf{C}\mathbf{M}^{-1}(\mathbf{M}\mathbf{x}_t + \mathbf{g}) \\ \mathbf{M}\mathbf{x}_{t+1} + \mathbf{g}|\mathbf{x}_t &\sim N_{Jp}(\mathbf{M}(\mathbf{A}\mathbf{x}_t + \mathbf{b}) + \mathbf{g}, \mathbf{M}\mathbf{Q}\mathbf{M}') \\ &= N_{Jp}(\mathbf{M}\mathbf{A}\mathbf{M}^{-1}(\mathbf{M}\mathbf{x}_t + \mathbf{g}) + (\mathbf{M}\mathbf{b} + \mathbf{g} - \mathbf{M}\mathbf{A}\mathbf{M}^{-1}\mathbf{g}), \mathbf{M}\mathbf{Q}\mathbf{M}')\end{aligned}$$

That means, when  $\{\mathbf{d}, \mathbf{C}, \mathbf{x}_t, \mathbf{A}, \mathbf{b}, \mathbf{Q}\}$  is a solution,  $\{\mathbf{d}^*, \mathbf{C}^*, \mathbf{x}_t^*, \mathbf{A}^*, \mathbf{b}^*, \mathbf{Q}^*\}$  is also a solution, where  $\mathbf{d}^* = \mathbf{d} - \mathbf{C}\mathbf{M}^{-1}\mathbf{g}$ ,  $\mathbf{C}^* = \mathbf{C}\mathbf{M}^{-1}$ ,  $\mathbf{x}_t^* = \mathbf{M}\mathbf{x}_t + \mathbf{g}$ ,  $\mathbf{A}^* = \mathbf{M}\mathbf{A}\mathbf{M}^{-1}$ ,  $\mathbf{b}^* = \mathbf{M}\mathbf{b} + \mathbf{g} - \mathbf{M}\mathbf{A}\mathbf{M}^{-1}\mathbf{g}$  and  $\mathbf{Q}^* = \mathbf{M}\mathbf{Q}\mathbf{M}'$ , for any block diagonal  $\mathbf{M}$ . Therefore, we must add some constraints to ensure the unique solution/stationary posterior!

## 5.1 Constraints

An affine transformation contains translation and linear transformation.

### 5.1.1 Translation

There are two strategies: (1) just ignore  $\mathbf{d}$  and let  $\mathbf{C}\mathbf{x}_t$  absorb all the intercepts; (2) center  $\mathbf{x}_t$  such that mean for each row of  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$  is 0. However, because of the block-diagonal constraints on  $\mathbf{M}$ , the strategy (1) may be problematic. Moreover, when  $\mathbf{d}$  is non-zero, the  $\mathbf{x}_t$  will look ugly. Another benefit for (2) is that we can use  $\mathbf{d}$  to help clustering. However, the constraint on translation cannot resolve problems from linear transformation, since any linear transformation will preserve the zero row mean in  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$ .

### 5.1.2 Linear Transformation

To mimic PCA, one natural way is to make the loading orthogonal. Again, I'm

considering two strategies: (1) make  $\mathbf{C}'^{(j)}\mathbf{C}^{(j)} = \mathbf{I}_p$ ; (2) make  $(\mathbf{C}'^{(1)}, \dots, \mathbf{C}'^{(J)}) \begin{pmatrix} \mathbf{C}^{(1)} \\ \vdots \\ \mathbf{C}^{(J)} \end{pmatrix} =$

$\mathbf{I}_p$ . Personally, I prefer the first one, for both computational and clustering reasons. Computationally, if I do things cluster-wise, I can easily use parallel when dimension becomes large.

**However, this is not enough!** Consider  $\mathbf{x}_t^* = \mathbf{U}\mathbf{x}_t$ , to preserve the orthogonality and block diagonal on  $\mathbf{C}$ ,  $\mathbf{U}$  must be block diagonal and each diagonal block  $\mathbf{U}^{(j)}$  is an orthogonal matrix (need to be the same for constraint (2)). More explicitly,

$$\begin{pmatrix} \mathbf{C}^{(1)} & \dots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \dots & \mathbf{C}^{(J)} \end{pmatrix} \begin{pmatrix} \mathbf{U}'^{(1)} & \dots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \dots & \mathbf{U}'^{(J)} \end{pmatrix} = \begin{pmatrix} \mathbf{C}^{(1)}\mathbf{U}'^{(J)} & \dots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \dots & \mathbf{C}^{(1)}\mathbf{U}'^{(J)} \end{pmatrix}$$

This transformation will also preserve the Frobenius norm of  $\mathbf{x}_t$ , but the element value is not preserved. To ensure the unique solution, we need to add more constraints on linear dynamics, i.e.  $\mathbf{A}$  or/and  $\mathbf{Q}$ . Simply making  $\mathbf{A}$  or  $\mathbf{Q}$  be block diagonal or even diagonal with same absolute value (e.g.,  $\mathbf{I}_{Jp}$ ) is still not enough. One solution is to force at least one of  $\mathbf{A}$  and  $\mathbf{Q}$  be diagonal.

Take  $\mathbf{A}$  for example. Assume the true  $\mathbf{A}$  is diagonal and true  $\mathbf{C}^{(j)}$  has orthogonal column. As shown previous, the linear transformation  $\mathbf{U}$  must be block diagonal and each diagonal block  $\mathbf{U}^{(j)}$  is an orthogonal matrix. Further, to preserve the diagonal structure of  $\mathbf{A}$  (the elements on the diagonal are generally not the same), we need  $\mathbf{U}^{(j)}$  to be diagonal. These two constraints force  $\mathbf{U}$  be diagonal with diagonal elements be  $\pm 1$ . So the solution is unique up to sign. The same rationale for constraint on  $\mathbf{Q}$ .

**The above argument reminds me that I may need one more tiny constraint to solve the sign flipping issue. Do it later.**

In summary, we need add three constraints: (1) translation of  $\mathbf{x}_t$ , (2) orthogonality of loading and (3) diagonality of linear dynamics. In the following implementation, I tried 3 sets of constraints:

- (1) row mean in  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$  is 0,  $\mathbf{C}'^{(j)}\mathbf{C}^{(j)} = \mathbf{I}_p$ ,  $\mathbf{A}$  is full while  $\mathbf{Q}$  is diagonal
- (2) row mean in  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$  is 0,  $\mathbf{C}'^{(j)}\mathbf{C}^{(j)} = \mathbf{I}_p$ ,  $\mathbf{Q}$  is full while  $\mathbf{A}$  is diagonal
- (3) row mean in  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$  is 0,  $\mathbf{C}'^{(j)}\mathbf{C}^{(j)} = \mathbf{I}_p$ , both  $\mathbf{A}$  and  $\mathbf{Q}$  are diagonal

Constraint (1) and (2) allows interaction between clusters.

Although, constraint (1) (diagonal  $\mathbf{Q}$ ) is valid in theory, it doesn't work well in practice. Since elements  $\mathbf{Q}$  is usually very small ( $10^{-3}$  to  $10^{-5}$ ). Therefore, I will discard this constraint later.

## 5.2 Constraint Implementation

### 5.2.1 Translation

The constraint on translation is easy, just row-center the sample of  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$  in each iteration. We can get the transformed full conditional distribution after

some algebra. Let  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T) = \begin{pmatrix} \mathbf{X}'_{(1)} \\ \vdots \\ \mathbf{X}'_{(J_p)} \end{pmatrix}$ , where  $\mathbf{x}_{(k)} \in \mathbb{R}^T$  is the k-th

row vector. Then the k-th row mean is  $\bar{\mathbf{x}}_{(k)} = \frac{1}{T}\mathbf{x}'_{(k)}\mathbf{1}_T$ . The unconstraint full conditional distribution of  $vec(\mathbf{X})$  is derived in previous section. After normalization in each iteration, we are actually sampling from  $P(vec(\mathbf{X}) -$

$\mathbf{1}_T \otimes \begin{pmatrix} \bar{\mathbf{x}}_{(1)} \\ \vdots \\ \bar{\mathbf{x}}_{(J_p)} \end{pmatrix} | \mathbf{Y}, \dots)$ . Notice that,

$$\mathbf{1}_T \otimes \begin{pmatrix} \bar{\mathbf{x}}_{(1)} \\ \vdots \\ \bar{\mathbf{x}}_{(J_p)} \end{pmatrix} = \mathbf{1}_T \otimes \frac{1}{T}\mathbf{X}\mathbf{1}_T = \frac{1}{T}vec\left(\mathbf{1}\mathbf{X}_T\mathbf{1}'_T\right) = \frac{\mathbf{J}_T \otimes \mathbf{I}_{J_p}}{T}vec(\mathbf{X})$$

, then we are just sampling from full conditional distribution of  $(\mathbf{I}_{J_p T} - \frac{\mathbf{J}_T \otimes \mathbf{I}_{J_p}}{T})vec(\mathbf{X})$

### 5.2.2 Linear Transformation

The constraint on orthogonality is kind of tricky. Here, I just transformed  $\mathbf{C}^{(j)}$  be  $\mathbf{C}^{(j)} = \mathbf{K}^{(j)} \left( \mathbf{K}'^{(j)} \mathbf{K}^{(j)} \right)^{-1/2}$ . Therefore, we can just sample  $\mathbf{K}^{(j)}$  without any constrained. Since  $\mathbf{K}^{(j)}$  has the same structure as  $\mathbf{C}^{(j)}$ , we can put some hierarchical prior to model  $\mathbf{C}^{(j)}$ . Specially, when  $\mathbf{K}^{(j)} \sim MN_{n,p}(\mathbf{0}, \Sigma^{(j)}, \mathbf{I})$  ( $MN_{n,p}$  means a matrix normal distribution),  $\mathbf{C}^{(j)} \sim MACG(\Sigma^{(j)})$ .  $MACG(\Sigma)$  means the matrix angular central Gaussian distribution, with density  $f_{\mathbf{X}}(\mathbf{X}) = |\Sigma|^{-p/2} |\mathbf{X}' \Sigma^{-1} \mathbf{X}|^{-n/2}$ . When  $\Sigma^{(j)} = \mathbf{I}_n$  or  $n = p$ ,  $\mathbf{C}^{(j)}$  is just uniformly distributed on the Stiefel manifold  $S(p, n) = \{\mathbf{C} \in \mathbb{R}^{n \times p} | \mathbf{C}' \mathbf{C} = \mathbf{I}_p\}$ .

**This reminds me to fix the problem when  $n < p$ , which may happen when**

doing clustering. One easy solution is just use the constraint  $(\mathbf{C}'^{(1)}, \dots, \mathbf{C}'^{(J)}) \begin{pmatrix} \mathbf{C}^{(1)} \\ \vdots \\ \mathbf{C}^{(J)} \end{pmatrix} =$

$\mathbf{I}_p$ , or just merge clusters with  $n < p$  to other clusters. Fix that.

Since now there's no closed form for full conditional, I just sample it by no U-turn sampler (NUTS).

The diagonality in linear dynamics is trivial, except for the constraint with diagonal  $\mathbf{A}$  and full  $\mathbf{Q}$ . This needs some algebra. Denote the diagonal  $\mathbf{A}$  as  $\mathbf{A} = \text{diag}(a_1, \dots, a_{Jp})$ . Then,

$$\mathbf{Y}_{BA} = \begin{pmatrix} \mathbf{x}_2' \\ \mathbf{x}_3' \\ \vdots \\ \mathbf{x}_T' \end{pmatrix} = \begin{pmatrix} 1 & x_{1,1} & \cdots & x_{1,Jp} \\ 1 & x_{2,1} & \cdots & x_{2,Jp} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{T-1,1} & \cdots & x_{T-1,Jp} \end{pmatrix} \begin{pmatrix} b_1 & \cdots & b_{Jp} \\ a_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & a_{Jp} \end{pmatrix} + \mathbf{E} = \mathbf{X}_{BA} \mathbf{F}_{BA} + \mathbf{E}$$

, where  $\mathbf{E} = (\epsilon_1, \dots, \epsilon_{T-1})'$  and  $\epsilon_t \sim N(\mathbf{0}, \mathbf{Q})$ . Notice that  $\mathbf{x}_t = (x_{t,1}, \dots, x_{t,Jp})'$ .

Equivalently, we can rewrite the model as:

$$\mathbf{Y}_{ba} = \text{vec}(\mathbf{Y}_{BA}) = \begin{pmatrix} 1 & x_{1,1} & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{T-1,1} & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & x_{1,Jp} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & x_{T-1,Jp} \end{pmatrix} \begin{pmatrix} b_1 \\ a_1 \\ \vdots \\ b_{Jp} \\ a_{Jp} \end{pmatrix} + \text{vec}(\mathbf{E}) = \mathbf{X}_{ba} \mathbf{f}_{ba} + \text{vec}(\mathbf{E})$$

, where  $\text{Cov}(\text{vec}(\mathbf{E})) = \mathbf{Q} \otimes \mathbf{I}_{T-1}$ .

If we use the prior:

$$\begin{aligned}\mathbf{Q} &\sim W^{-1}(\Psi_{Q_0}, \nu_{Q_0}) \\ \mathbf{f}_{ba}|\mathbf{Q} &\sim N_{2Jp}(\mathbf{1}_{Jp} \otimes \mathbf{f}_{ba,0}, \mathbf{Q} \otimes \Lambda_{ba,0}^{-1})\end{aligned}$$

, where  $\mathbf{f}_{ba,0} = (0, 1)'$  and  $\Lambda_{ba,0} = \mathbf{I}_2$ . Then the corresponding prior for  $\text{vec}(\mathbf{F}_{BA})$  is:

$$\text{vec}(\mathbf{F}_{BA})|\mathbf{Q} \sim N(\text{vec}(\mathbf{F}_{BA,0}), \mathbf{Q} \otimes \Lambda_{BA,0}^{-1})$$

, where  $\mathbf{F}_{BA,0} = \begin{pmatrix} \mathbf{0}_{Jp}' \\ \mathbf{I}_{Jp} \end{pmatrix}$  and  $\Lambda_{BA,0} = \mathbf{I}_{Jp+1}$ .

Then the posteriors are:

$$\begin{aligned}\mathbf{Q}|\mathbf{Y}, \dots &\sim W^{-1}(\Psi_Q, \nu_Q) \\ \mathbf{f}_{ba}|\mathbf{Y}, \mathbf{Q}, \dots &\sim N(\hat{\mathbf{f}}_{ba}, (\mathbf{Q} \otimes \mathbf{I}_2) \Lambda_{ba}^{-1})\end{aligned}$$

, with parameters be:

$$\begin{aligned}\Lambda_{ba} &= \mathbf{X}_{ba}' \mathbf{X}_{ba} + \mathbf{I}_{Jp} \otimes \Lambda_{ba,0} \\ \hat{\mathbf{f}}_{ba} &= \Lambda_{ba}^{-1} (\mathbf{X}_{ba}' \mathbf{Y}_{ba} + \mathbf{1}_{Jp} \otimes \Lambda_{ba,0} \mathbf{f}_{ba,0})\end{aligned}$$

The corresponding estimate for  $\mathbf{F}_{BA}$  is denoted as  $\hat{\mathbf{F}}_{BA}$ , which is just a transformation of  $\hat{\mathbf{f}}_{ba}$ .

$$\begin{aligned}\Psi_Q &= \Psi_{Q_0} + \left( \mathbf{Y}_{BA} - \mathbf{X}_{BA} \hat{\mathbf{F}}_{BA} \right)' \left( \mathbf{Y}_{BA} - \mathbf{X}_{BA} \hat{\mathbf{F}}_{BA} \right) + \left( \hat{\mathbf{F}}_{BA} - \mathbf{F}_{BA,0} \right)' \Lambda_{BA,0}^{-1} \left( \hat{\mathbf{F}}_{BA} - \mathbf{F}_{BA,0} \right) \\ \nu_Q &= \nu_{Q_0} + T - 1\end{aligned}$$

### 5.3 Simulation

The simulation example is as before, and all the code can be found in this folder.

Since now we constrain the loading has orthogonal column for each cluster, just showing the Frobenius norm of  $\mathbf{x}_t$  is not enough to check convergence. Here, I further check the value of  $\mathbf{x}_{[T/2]}$ .

For  $\mathbf{C}^{(j)} = \mathbf{K}^{(j)} \left( \mathbf{K}'^{(j)} \mathbf{K}^{(j)} \right)^{-1/2}$ , I just put a simple hierarchical prior on  $\mathbf{K}^{(j)}$  such that each element of  $\mathbf{K}^{(j)}$  is i.i.d.  $N(\mu_k^{(j)}, (\sigma_k^{(j)})^2)$ . Well, maybe we can put some more careful priors later to reflect the covariance among neurons.

Basically, I ran both MCMCs for 10,000 iterations. Here, I show 6 sets of plots:

- (1) Full trace of Frobenius norm in  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T)$  and the value of  $\mathbf{x}_{T/2} = \mathbf{x}_{500}$
- (2) Full trace of 2-norm/ Frobenius norm in  $\mathbf{d}$ ,  $\mathbf{b}$ ,  $\mathbf{A}$  and  $\mathbf{Q}$
- (3) partial trace plot of (1) from iteration 1 to 100
- (4) partial trace plot of (2) from iteration 1 to 100
- (5) average mean firing rate from iteration 5000 to 10,000.
- (6) average  $\mathbf{X}$  from iteration 5000 to 10,000.

### 5.3.1 $\mathbf{Q}$ is full while $\mathbf{A}$ is diagonal

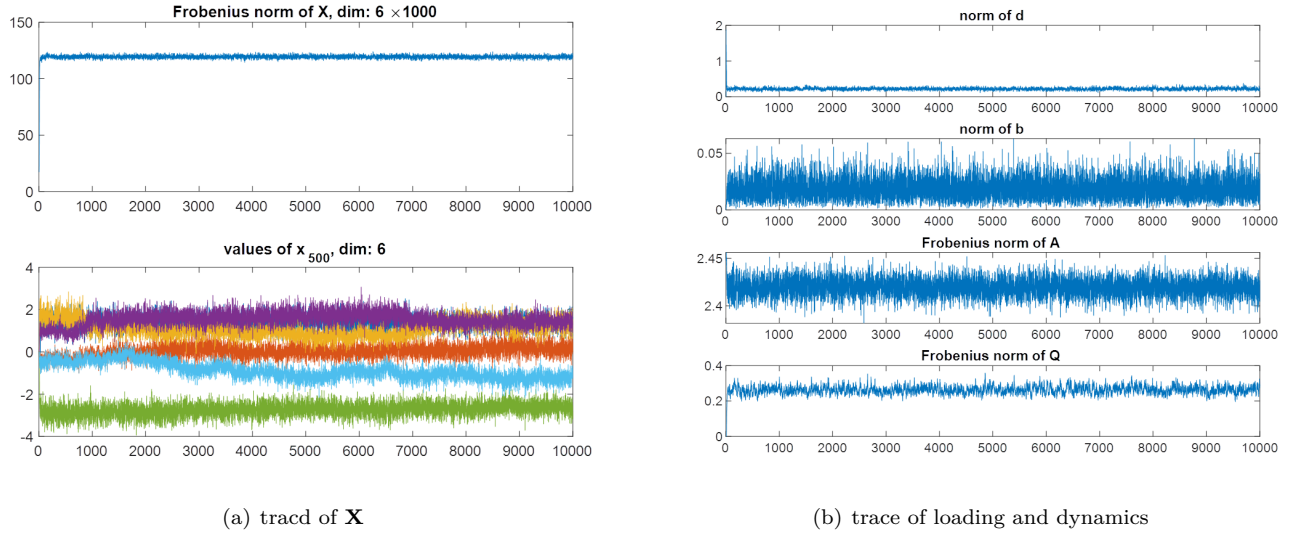
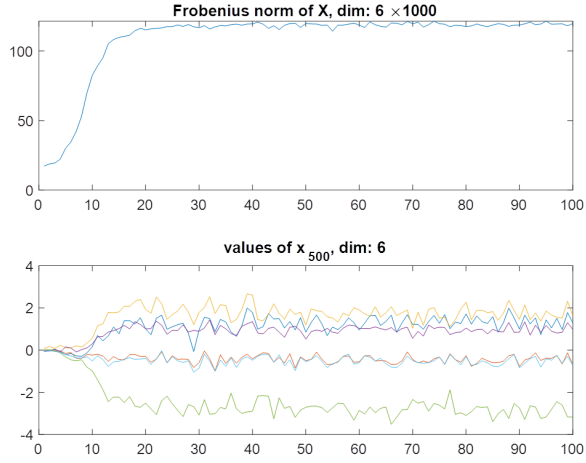
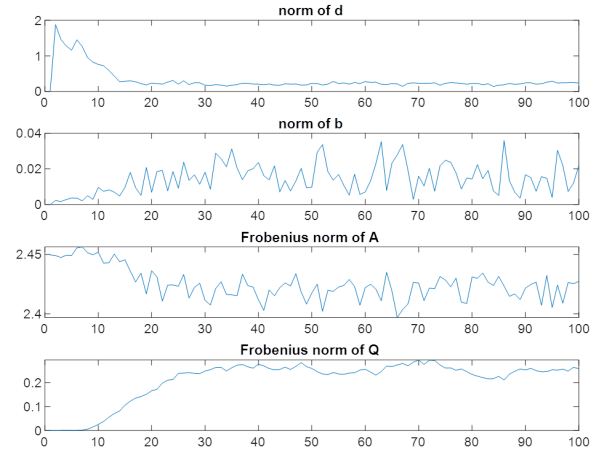


Figure 1:  $\mathbf{Q}$  is full while  $\mathbf{A}$  is diagonal, full trace

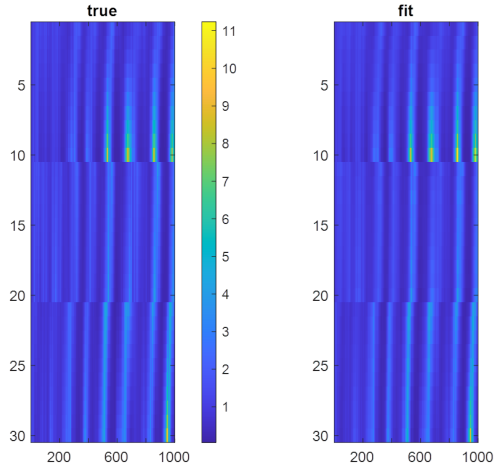


(a) trace of  $\mathbf{X}$

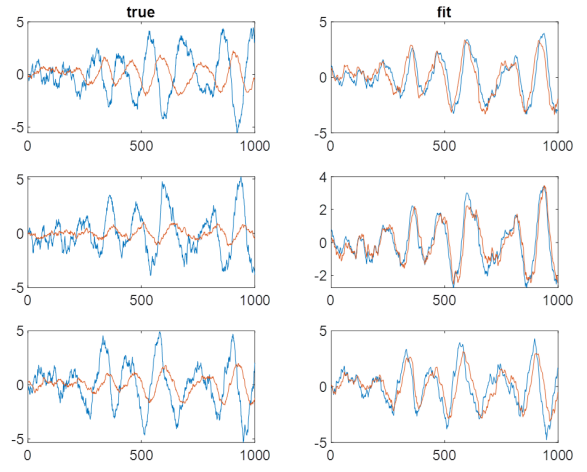


(b) trace of loading and dynamics

Figure 2:  $\mathbf{Q}$  is full while  $\mathbf{A}$  is diagonal, partial trace



(a) mean firing rate

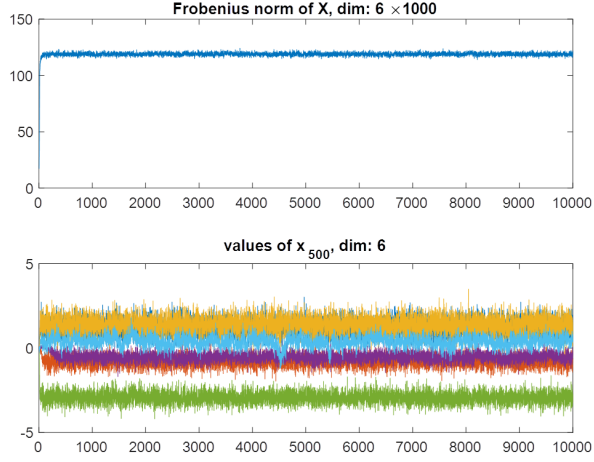


(b) latent vectors

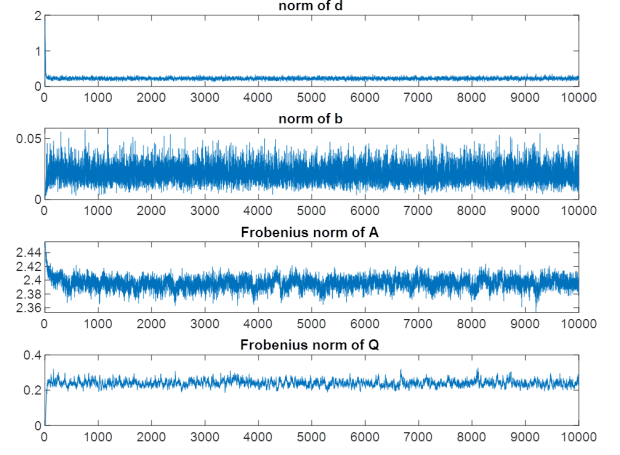
Figure 3:  $\mathbf{Q}$  is full while  $\mathbf{A}$  is diagonal, mean FR and  $\mathbf{x}_t$



### 5.3.2 both $\mathbf{A}$ and $\mathbf{Q}$ are diagonal

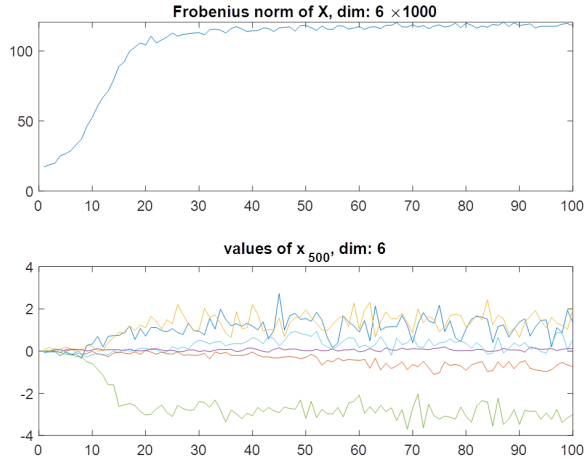


(a) trace of  $\mathbf{X}$

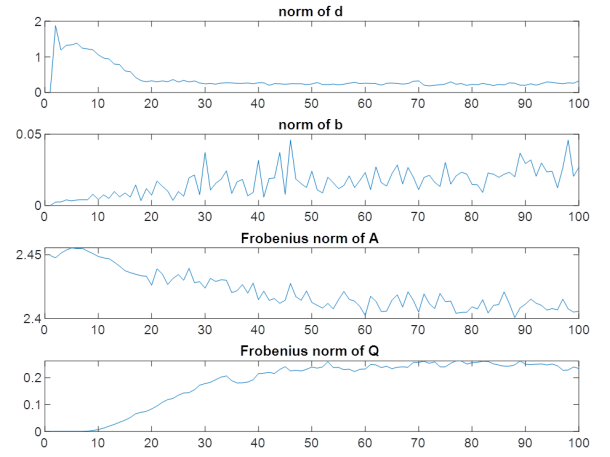


(b) trace of loading and dynamics

Figure 4: both  $\mathbf{A}$  and  $\mathbf{Q}$  are diagonal, full trace



(a) trace of  $\mathbf{X}$



(b) trace of loading and dynamics

Figure 5: both  $\mathbf{A}$  and  $\mathbf{Q}$  are diagonal, partial trace

We can see that either is fine, and the convergence is achieved at nearly 20 to 30 iterations.

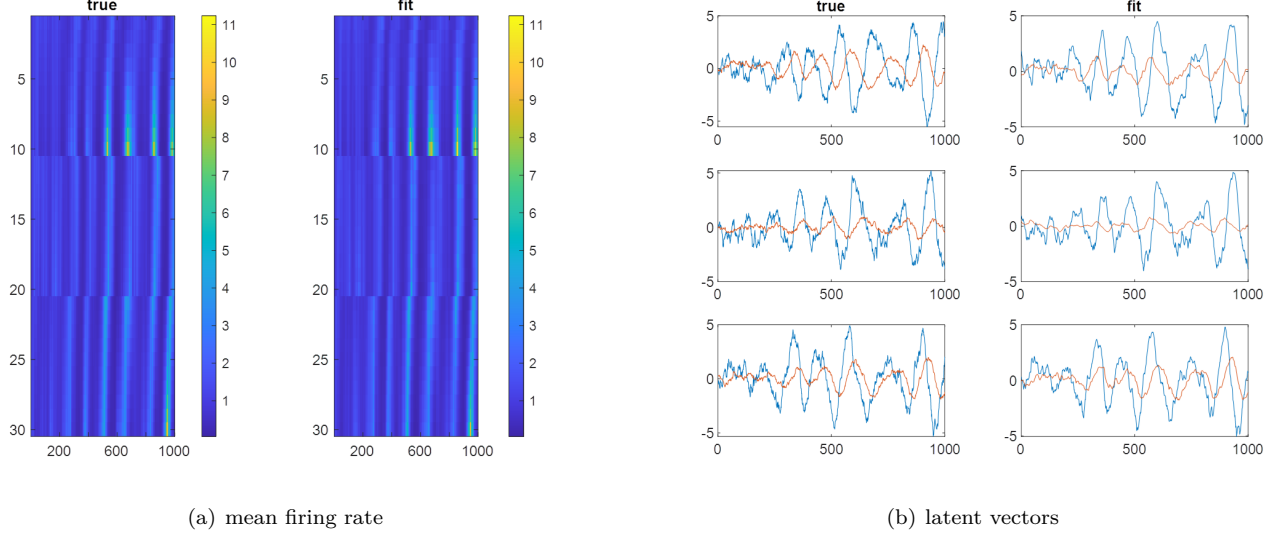


Figure 6: both  $\mathbf{A}$  and  $\mathbf{Q}$  are diagonal, mean FR and  $\mathbf{x}_t$

**In summary, to ensure convergence, we can put the following two constraints:**

- (1) row mean in  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$  is 0,  $\mathbf{C}'^{(j)}\mathbf{C}^{(j)} = \mathbf{I}_p$ ,  $\mathbf{Q}$  is full while  $\mathbf{A}$  is diagonal
- (2) row mean in  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$  is 0,  $\mathbf{C}'^{(j)}\mathbf{C}^{(j)} = \mathbf{I}_p$ , both  $\mathbf{A}$  and  $\mathbf{Q}$  are diagonal

If we really care about the interactions between clusters, just use constraint (1). Otherwise, either is fine.

#### 5.4 Comments on Previous Models

When fitting the LDS or factor analysis (e.g. GPFA) model, I may only believe the temporal pattern or the change of the pattern. The pattern itself at a specific time point doesn't make a lot of sense, because of the affine transformation invariance.

In previous work, they fit models by EM, variational Bayes or some other MAP/ MLE method. **This is quite dangerous** if they don't check the uniqueness (Thanks to the gold standard MCMC, otherwise I will never check this issue...).

In the PLDS code that Ian shared me before, the default model has no constraint. They fit it by EM, and use ELBO as the convergence criteria. Although the parameters are not unique, the likelihood keeps the same, using ELBO is justified. However, since there's no randomness in EM, it's impossible to detect the non-uniqueness. Actually, I think their results are even not local optima,

but just a single point in the optimal surface (I guess if they use something like stochastic EM, they will find the issue).

The PLDS also provides some constraint options, but some of them are still not enough to ensure unique solution. BTW, this might be another reason they choose to use EM, since adding constraint in EM is trivial: just do unconstrained EM in each iteration and project the value to the constraint space (i.e. add constraint) at the end of each iteration. But this doesn't work for MCMC, since we are sampling the distribution but not just the single optimal value.

**I guess similar problems may happen in other LDS/ GPFA research.**

Even if they make enough constraints to ensure uniqueness, these research should be careful about the pattern they found. For example, in Joshua et al., 2020, the changes/ different states of interaction make sense to me, but I'm quite skeptical about their results in the dynamic matrix  $\mathbf{A}$  itself. Maybe because of block-diagonal structure in loading  $\mathbf{C}$ , the interaction between population make sense somewhat. But at least the relationship within population doesn't make any sense, because of transformation invariance.

## 6 TODO

- (1) Improve the NUTS.
- (2) Debug the clustering.
- (3) Check and clear typo.
- (4) Implement the mixture of finite mixture (MFM) by Jeff Miller.
- (5) Find a more efficient way to generate new cluster parameters, otherwise the newly generated ones will always be rejected.
- (6) After all of them are resolved, switch to GP version