

# Complex Synapse Identification

Subhaneil Lahiri

August 4, 2020

## Abstract

Notes on using HMM techniques to fit models of complex synapses

## 1 Introduction

## 2 Notation

Let  $t$  be the time of a plasticity event, assumed to be an integer. We denote a sequence of event times by  $[0, T]$ . Any sequence whose upper bound is smaller than its lower bound can be considered empty, and that variable is omitted. The same applies to sums and products.

We denote states of a synapse by subscripts  $i, j, k$ . The state at a particular time during a sequence is  $i(t)$ . An entire trajectory (sequence of states) is denoted by  $i[0, T]$ .

The synaptic weight is denoted by  $\mathbf{w}$ . Its value when the synapse is in state  $i$  is  $\mathbf{w}_i$ . Its observed value at time  $t$  during a sequence is  $\mathbf{w}(t)$ . An entire sequence of observations is denoted by  $\mathbf{w}[0, T]$ .

The different types of plasticity (e.g. potentiation and depression) are indicated by superscripts  $\mu, \nu$ . The type of plasticity used between times  $(t, t + 1)$  is  $\mu(t)$ . An entire sequence of plasticity types is denoted by  $\mu[0, T - 1]$ .

When we are combining data from several sequences of plasticity events, the individual sequences will be denoted by superscripts  $a, b$ . The set of all sequences will be denoted by the superscript  $A$ .

The Markov matrices describing plasticity are  $\mathbf{M}_{ij}^\mu$ . The probabilities of the initial states are  $\boldsymbol{\pi}_i$ . The set of parameters defining a model (i.e.  $\{\mathbf{M}_{ij}^\mu\}$  and  $\boldsymbol{\pi}_i$ ) are denoted by  $\mathbf{M}$  for convenience.

sec:em

### 3 Expectation-Maximisation (EM)

When we don't have any prior information about the model  $\mathbf{M}$ , we fit the model by maximising the likelihood function:

$$\begin{aligned}\mathbb{P}_{[0,T]}^A(\mathbf{w} \mid \mu, \mathbf{M}) &\equiv \mathbb{P}(\mathbf{w}^A[0, T] \mid \mu^A[0, T-1], \mathbf{M}) \\ &= \prod_a \mathbb{P}(\mathbf{w}^a[0, T] \mid \mu^a[0, T-1], \mathbf{M}) \equiv \prod_a \mathbb{P}_{[0,T]}^a(\mathbf{w} \mid \mu, \mathbf{M}).\end{aligned}\tag{1} \quad \text{\texttt{eq:like}}$$

To do this, we define an axillary function [1, 2]:

$$Q(\mathbf{M}, \mathbf{M}') = \sum_{i^A[0,T]} \mathbb{P}_{[0,T]}^A(i, \mathbf{w} \mid \mu, \mathbf{M}) \log \left[ \frac{\mathbb{P}_{[0,T]}^A(i, \mathbf{w} \mid \mu, \mathbf{M}')}{\mathbb{P}_{[0,T]}^A(i, \mathbf{w} \mid \mu, \mathbf{M})} \right].\tag{2} \quad \text{\texttt{eq:emq}}$$

It can be shown that, if  $\mathbf{M}'$  is chosen to maximise this quantity (which is easier than maximising the likelihood), then it will have a higher likelihood than  $\mathbf{M}$ . In effect,  $\mathbf{M}$  is used to estimate the hidden state trajectories, which are then used to re-estimate the model. This can be done iteratively to find a local maximum of the likelihood. If this is repeated several times with random initialisation, we can hope to find the global maximum.

When we have several sequences, we can write

$$\begin{aligned}Q(\mathbf{M}, \mathbf{M}') &= \sum_{i^A[0,T]} \prod_b \mathbb{P}_{[0,T]}^b(i, \mathbf{w} \mid \mu, \mathbf{M}) \sum_a \log \left[ \frac{\mathbb{P}_{[0,T]}^a(i, \mathbf{w} \mid \mu, \mathbf{M}')}{\mathbb{P}_{[0,T]}^a(i, \mathbf{w} \mid \mu, \mathbf{M})} \right] \\ &= \sum_a \prod_{b \neq a} \mathbb{P}_{[0,T]}^b(\mathbf{w} \mid \mu, \mathbf{M}) \sum_{i^a[0,T]} \mathbb{P}_{[0,T]}^a(i, \mathbf{w} \mid \mu, \mathbf{M}) \log \left[ \frac{\mathbb{P}_{[0,T]}^a(i, \mathbf{w} \mid \mu, \mathbf{M}')}{\mathbb{P}_{[0,T]}^a(i, \mathbf{w} \mid \mu, \mathbf{M})} \right] \\ &= \mathbb{P}_{[0,T]}^A(\mathbf{w} \mid \mu, \mathbf{M}) \sum_a \frac{Q^a(\mathbf{M}, \mathbf{M}')}{\mathbb{P}_{[0,T]}^a(\mathbf{w} \mid \mu, \mathbf{M})}.\end{aligned}\tag{3} \quad \text{\texttt{eq:emqsever}}$$

sec:bw

#### 3.1 Single sequence (Baum-Welch)

Suppose we have a single sequence of plasticity events. Then the joint likelihood of a trajectory and sequence observed synaptic weights is given by

$$\mathbb{P}(i[0, T], \mathbf{w}[0, T] \mid \mu[0, T-1], \mathbf{M}) = \pi_{i(0)} \Pi_{i(0)}(0) \prod_{t=0}^{T-1} \mathbf{M}_{i(t)i(t+1)}^{\mu(t)} \Pi_{i(t+1)}(t+1),\tag{4} \quad \text{\texttt{eq:trajlike}}$$

where  $\Pi_i(t) = \delta_{\mathbf{w}_i \mathbf{w}(t)}$

We wish to maximise (2) subject to normalisation constraints. Therefore we maximise the Lagrangian:

$$\mathcal{L} = Q(\mathbf{M}, \mathbf{M}') + \lambda \left( 1 - \sum_i \pi'_i \right) + \sum_{\mu i} \lambda_i^\mu \left( 1 - \sum_j \mathbf{M}'_{ij}^\mu \right).\tag{5} \quad \text{\texttt{eq:lagrange}}$$

We find

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{M}_{ij}^{\mu}} &= \frac{N_{ij}^{\mu}}{\mathbf{M}_{ij}^{\mu}} - \lambda_i^{\mu}, \\ \frac{\partial \mathcal{L}}{\partial \boldsymbol{\pi}'_i} &= \frac{\alpha_i(0) \beta_i(0)}{\boldsymbol{\pi}'_i} - \lambda,\end{aligned}\tag{6}$$

where the numerator  $N_{ij}^{\mu}$  and the forward/backward variables  $\alpha$  and  $\beta$  are defined as:

$$\begin{aligned}N_{jk}^{\nu} &= \sum_{t=0}^{T-1} \delta_{\mu(t)\nu} \mathbb{P}(i[t, t+1] = (j, k), \mathbf{w}[0, T] \mid \mu[0, T-1], \mathbf{M}) \\ &= \sum_{t=1}^T \delta_{\mu(t)\nu} \alpha_j(t) \beta_k(t+1) \mathbf{M}_{jk}^{\nu} \boldsymbol{\Pi}_k(t+1), \\ \alpha_j(t) &= \mathbb{P}(i(t) = j, \mathbf{w}[0, t] \mid \mu[0, t-1], \mathbf{M}) \\ &= \sum_{i[0, t]} \boldsymbol{\pi}_{i(0)}(0) \boldsymbol{\Pi}_{i(0)}(0) \prod_{s=0}^{t-1} \mathbf{M}_{i(s)i(s+1)}^{\mu(s)} \boldsymbol{\Pi}_{i(s+1)}(s+1) \delta_{i(t)j}, \\ \beta_j(t) &= \mathbb{P}(\mathbf{w}[t+1, T] \mid i(t) = j, \mu[t, T-1], \mathbf{M}) \\ &= \sum_{i[t, T]} \delta_{i(t)j} \prod_{s=t}^{T-1} \mathbf{M}_{i(s)i(s+1)}^{\mu(s)} \boldsymbol{\Pi}_{i(s+1)}(s+1).\end{aligned}\tag{7}$$

Note that, whilst disjoint sequences  $\mathbf{w}[s, s']$  and  $\mathbf{w}[t, t']$  are not independent, they are conditionally independent given the state at any time between the two intervals. Therefore:

$$\alpha_j(t) \beta_j(t) = \mathbb{P}(i(t) = j, \mathbf{w}[0, T] \mid \mathbf{M}).\tag{8}$$

These quantities usually suffer from underflow. This can be avoided by using normalised forward/backward variables [3]:

$$\begin{aligned}\tilde{\alpha}_i(t) &= \alpha_i(t) \prod_{s=0}^t \eta(s), \\ \tilde{\beta}_i(t) &= \beta_i(t) \prod_{s=t+1}^T \eta(s),\end{aligned}\tag{9}$$

where  $\eta[0, T]$  is chosen so that all of the  $\tilde{\alpha}[0, T]$  are normalised:  $\sum_i \tilde{\alpha}_i(t) = 1$ . These

quantities have the interpretations:

$$\begin{aligned}
\tilde{\alpha}_j(t) &= \mathbb{P}(i(t) = j \mid \mathbf{w}[0, t], \mu[0, t-1], \mathbf{M}), \\
\eta(t)^{-1} &= \mathbb{P}(\mathbf{w}(t) \mid \mathbf{w}[0, t-1], \mu[0, t-1], \mathbf{M}), \\
\left[ \prod_{s=0}^t \eta(s) \right]^{-1} &= \mathbb{P}(\mathbf{w}[0, t] \mid \mu[0, t-1], \mathbf{M}), \\
\left[ \prod_{s=t_0}^{t_1} \eta(s) \right]^{-1} &= \mathbb{P}(\mathbf{w}[t_0, t_1] \mid \mathbf{w}[0, t_0-1], \mu[0, t_1-1], \mathbf{M}).
\end{aligned} \tag{10} \quad \text{\texttt{\{eq:norminte}}}$$

Note that the new backward variables  $\tilde{\beta}$  are not normalised and have no particularly interesting interpretation, beyond the following:

$$\tilde{\alpha}_j(t) \tilde{\beta}_j(t) = \mathbb{P}(i(t) = j \mid \mathbf{w}[0, T], \mu[0, T-1], \mathbf{M}). \tag{11} \quad \text{\texttt{\{eq:normstat}}}$$

Then we can write

$$\begin{aligned}
\tilde{N}_{jk}^\nu &= \sum_{t=0}^{T-1} \delta_{\mu(t)\nu} \mathbb{P}(i[t, t+1] = (jk) \mid \mathbf{w}[0, T], \mu[0, T-1], \mathbf{M}) \\
&= \sum_{t=0}^{T-1} \delta_{\mu(t)\nu} \tilde{\alpha}_j(t) \tilde{\beta}_k(t+1) \eta(t+1) \mathbf{M}_{jk}^\nu \mathbf{\Pi}_k(t+1), \\
\frac{\partial \mathcal{L}}{\partial \mathbf{M}_{ij}'^\mu} &= \mathbb{P}_{[0, T]}(\mathbf{w} \mid \mu, \mathbf{M}) \frac{\tilde{N}_{ij}^\mu}{\mathbf{M}_{ij}'^\mu} - \lambda_i^\mu, \\
\frac{\partial \mathcal{L}}{\partial \boldsymbol{\pi}_i'} &= \mathbb{P}_{[0, T]}(\mathbf{w} \mid \mu, \mathbf{M}) \frac{\tilde{\alpha}_i(0) \tilde{\beta}_i(0)}{\boldsymbol{\pi}_i'} - \lambda,
\end{aligned} \tag{12} \quad \text{\texttt{\{eq:singledi}}}$$

The factors of  $\mathbb{P}_{[0, T]}(\mathbf{w} \mid \mu, \mathbf{M})$  can be absorbed by a redefinition of the Lagrange multipliers.

Demanding a maximum and setting the derivatives to zero leaves us with the update rule [1]:

$$\begin{aligned}
\mathbf{M}_{ij}'^\mu &= \frac{\tilde{N}_{ij}^\mu}{\sum_k \tilde{N}_{ik}^\mu}, \\
\boldsymbol{\pi}_i' &= \frac{\tilde{\alpha}_i(0) \tilde{\beta}_i(0)}{\sum_j \tilde{\alpha}_j(0) \tilde{\beta}_j(0)}.
\end{aligned} \tag{13} \quad \text{\texttt{\{eq:BWupdate}}}$$

This can be iterated until a maximum is found.

### 3.2 Multiple sequences (Rabiner-Juang)

When we have multiple sequences of plasticity events, we will still want to maximise the Lagrangian (5), but now we will use the auxiliary function (3). This results in

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \mathbf{M}_{ij}^{\prime\mu}} &= \mathbb{P}_{[0,T]}^A(\mathbf{w} \mid \mu, \mathbf{M}) \sum_a \frac{1}{\mathbb{P}_{[0,T]}^a(\mathbf{w} \mid \mu, \mathbf{M})} \frac{N_{ij}^{a\mu}}{\mathbf{M}_{ij}^{\prime\mu}} - \lambda_i^\mu, \\
&= \mathbb{P}_{[0,T]}^A(\mathbf{w} \mid \mu, \mathbf{M}) \frac{\sum_a \tilde{N}_{ij}^{a\mu}}{\mathbf{M}_{ij}^{\prime\mu}} - \lambda_i^\mu, \\
\frac{\partial \mathcal{L}}{\partial \pi_i'} &= \mathbb{P}_{[0,T]}^A(\mathbf{w} \mid \mu, \mathbf{M}) \sum_a \frac{1}{\mathbb{P}_{[0,T]}^a(\mathbf{w} \mid \mu, \mathbf{M})} \frac{\alpha_i^a(0) \beta_i^a(0)}{\pi_i'} - \lambda, \\
&= \mathbb{P}_{[0,T]}^A(\mathbf{w} \mid \mu, \mathbf{M}) \frac{\sum_a \tilde{\alpha}_i^a(0) \tilde{\beta}_i^a(0)}{\pi_i'} - \lambda.
\end{aligned} \tag{14}$$

{eq:multidif

Once again, we can absorb the factors of  $\mathbb{P}_{[0,T]}^A(\mathbf{w} \mid \mu, \mathbf{M})$  by a redefinition of the Lagrange multipliers.

Defining  $\tilde{N}_{ij}^{A\mu} = \sum_a \tilde{N}_{ij}^{a\mu}$  and setting the derivatives to zero leaves us with the update rule [4]:

$$\begin{aligned}
\mathbf{M}_{ij}^{\prime\mu} &= \frac{\tilde{N}_{ij}^{A\mu}}{\sum_k \tilde{N}_{ik}^{A\mu}}, \\
\pi_i' &= \frac{\sum_a \tilde{\alpha}_i^a(0) \tilde{\beta}_i^a(0)}{\sum_{bj} \tilde{\alpha}_j^b(0) \tilde{\beta}_j^b(0)}.
\end{aligned} \tag{15}$$

{eq:RJupdate

This can be iterated until a maximum is found.

## 4 Sparseness promoting priors

When we have a prior distribution for the model,  $\mathbb{P}(\mathbf{M})$ , instead of maximising the likelihood (1), we can maximise the posterior:

$$\mathbb{P}(\mathbf{M} \mid \mathbf{w}[0, T]) = \frac{\mathbb{P}(\mathbf{w}[0, T], \mathbf{M})}{\mathbb{P}(\mathbf{w}[0, T])} = \frac{\mathbb{P}(\mathbf{w}[0, T] \mid \mathbf{M}) \mathbb{P}(\mathbf{M})}{\mathbb{P}(\mathbf{w}[0, T])}. \tag{16}$$

{eq:posterior

As the denominator is independent of the model, this is equivalent to maximising the numerator, i.e. the joint distribution.

Using the same trick that was used for (2), we can define an auxilliary function

$$\begin{aligned}
\tilde{Q}(\mathbf{M}, \mathbf{M}') &= \sum_{i^A[0, T]} \mathbb{P}(i^A[0, T], \mathbf{w}^A[0, T], \mathbf{M}) \log \left[ \frac{\mathbb{P}(i^A[0, T], \mathbf{w}^A[0, T], \mathbf{M}')}{\mathbb{P}(i^A[0, T], \mathbf{w}^A[0, T], \mathbf{M})} \right] \\
&= \mathbb{P}(\mathbf{M}) Q(\mathbf{M}, \mathbf{M}') + \mathbb{P}(\mathbf{w}^A[0, T], \mathbf{M}) \log \left[ \frac{\mathbb{P}(\mathbf{M}')}{\mathbb{P}(\mathbf{M})} \right].
\end{aligned} \tag{17}$$

{eq:emqprior

We will consider priors of the form

$$\mathbb{P}(\mathbf{M}) = \frac{\exp\left(-\beta \sum_{\mu ij} E_{ij}^{\mu}(\mathbf{M}_{ij}^{\mu})\right)}{Z(\beta)}. \quad (18) \quad \{\text{eq:priorsum}\}$$

We will not put any prior on  $\boldsymbol{\pi}$ , so its update rule will be unaffected. The normalisation constant  $Z(\beta)$  will not play any role either. Now the derivatives read:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{M}_{ij}^{\mu}} = \mathbb{P}(\mathbf{w}[0, T], \mathbf{M}) \left( \frac{\tilde{N}_{ij}^{A\mu}}{\mathbf{M}_{ij}^{\mu}} - \beta \frac{\partial E_{ij}^{\mu}}{\partial \mathbf{M}_{ij}^{\mu}} \right) - \lambda_i^{\mu}. \quad (19) \quad \{\text{eq:priordif}\}$$

Once again, we can absorb the factors of  $\mathbb{P}(\mathbf{w}[0, T], \mathbf{M})$  by a redefinition of the Lagrange multipliers.

$\{\text{offdiagL1}\}$

## 4.1 Off-diagonal $\mathcal{L}^1$ penalty

Here we use the off-diagonal  $\mathcal{L}^1$  norm as a penalty:

$$E_{ij}^{\mu} = (1 - \delta_{ij}) \mathbf{M}_{ij}^{\mu}. \quad (20) \quad \{\text{eq:L1penalt}\}$$

Then the condition for a maximum (19) reads

$$\begin{aligned} \frac{\tilde{N}_{ij}^{A\mu}}{\mathbf{M}_{ij}^{\mu}} - (1 - \delta_{ij})\beta - \lambda_i^{\mu} &= 0, \\ \mathbf{M}_{ij}^{\mu} &= \frac{\tilde{N}_{ij}^{A\mu}}{\lambda_i^{\mu} + (1 - \delta_{ij})\beta}. \end{aligned} \quad (21) \quad \{\text{eq:L1priorm}\}$$

Demanding normalisation gives

$$\begin{aligned} \sum_{j \neq i} \frac{\tilde{N}_{ij}^{A\mu}}{\lambda_i^{\mu} + \beta} + \frac{\tilde{N}_{ii}^{A\mu}}{\lambda_i^{\mu}} &= 1, \\ \Rightarrow \quad \lambda_i^{\mu} &= \frac{1}{2} \left[ \left( \sum_j \tilde{N}_{ij}^{A\mu} - \beta \right) + \sqrt{\left( \sum_j \tilde{N}_{ij}^{A\mu} - \beta \right)^2 + 4\beta \tilde{N}_{ii}^{A\mu}} \right] \end{aligned} \quad (22) \quad \{\text{eq:L1norm}\}$$

Which produces the update rule

$$\begin{aligned} \mathbf{M}_{ij}^{\mu} &= \frac{2\tilde{N}_{ij}^{A\mu}}{\left( \sum_k \tilde{N}_{ik}^{A\mu} + \beta \right) + \sqrt{\left( \sum_k \tilde{N}_{ik}^{A\mu} - \beta \right)^2 + 4\beta \tilde{N}_{ii}^{A\mu}}} \quad \text{for } i \neq j, \\ \mathbf{M}_{ii}^{\mu} &= \frac{2\tilde{N}_{ii}^{A\mu}}{\left( \sum_k \tilde{N}_{ik}^{A\mu} - \beta \right) + \sqrt{\left( \sum_k \tilde{N}_{ik}^{A\mu} - \beta \right)^2 + 4\beta \tilde{N}_{ii}^{A\mu}}}. \end{aligned} \quad (23) \quad \{\text{eq:L1update}\}$$

sec:Lhalf

## 4.2 $\mathcal{L}^{\frac{1}{2}}$ penalty

Here we use the  $\mathcal{L}^{\frac{1}{2}}$  norm as a penalty:

$$E_{ij}^{\mu} = 2\sqrt{\mathbf{M}_{ij}^{\mu}}. \quad (24) \quad \text{eq:Lhalfpen}$$

Then the condition for a maximum (19) reads

$$\frac{\tilde{N}_{ij}^{A\mu}}{\mathbf{M}_{ij}^{\mu}} - \frac{\beta}{\sqrt{\mathbf{M}_{ij}^{\mu}}} - \lambda_i^{\mu} = 0. \quad (25) \quad \text{eq:Lhalfpri}$$

It is helpful to define

$$\begin{aligned} \tilde{N}_{ij}^{A\mu} &= \beta^2 L_{ij}^{\mu}, \\ \lambda_i^{\mu} &= (\gamma_i^{\mu})^{-2}. \end{aligned} \quad (26) \quad \text{eq:Lhalfvar}$$

Then we have the update rule

$$\mathbf{M}_{ij}^{\mu} = \frac{(\gamma_i^{\mu})^2 \beta^2}{4} \left( \sqrt{(\gamma_i^{\mu})^2 + 4L_{ij}^{\mu}} - \gamma_i^{\mu} \right)^2, \quad (27) \quad \text{eq:Lhalfupd}$$

with  $\gamma_i^{\mu}$  determined by:

$$\frac{(\gamma_i^{\mu})^2 \beta^2}{4} \sum_j \left( \sqrt{(\gamma_i^{\mu})^2 + 4L_{ij}^{\mu}} - \gamma_i^{\mu} \right)^2 = 1. \quad (28) \quad \text{eq:Lhalfnor}$$

## References

0baumwelch

[1] L. Baum, T. Petrie, G. Soules, and N. Weiss, “A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains,” *The annals of mathematical statistics* **41** (1970) no. 1, 164–171.

ster2007EM

[2] A. P. Dempster, N. M. Laird, and D. B. Rubin, “Maximum Likelihood from Incomplete Data via the EM Algorithm,” *Journal of the Royal Statistical Society. Series B (Methodological)* (Oct., 2007).

003HMMnorm

[3] C. Zhai, “A Brief Note on the Hidden Markov Models (HMMs).” Lecture notes, 2003. <http://sifaka.cs.uiuc.edu/course/498cxz05f/hmm.pdf>.

3speechrec

[4] L. R. Rabiner and B.-H. Juang, *Fundamentals of speech recognition*, vol. 14 of *Signal Processing*. Prentice Hall, Inc., Upper Saddle River, NJ, USA, 1993.