

# Research Notebook

Taisuke Yasuda

June 16, 2017

## Contents

<b>1</b>	<b>May 20, 2017</b>	<b>4</b>
1.1	Progress . . . . .	4
1.2	Statistical Model . . . . .	4
1.2.1	Justification . . . . .	4
<b>2</b>	<b>May 22, 2017</b>	<b>4</b>
2.1	Progress . . . . .	4
2.2	Point Estimation of Model Parameters . . . . .	4
2.2.1	Estimations Under a Simplified Model . . . . .	4
2.2.2	EM Algorithm . . . . .	5
<b>3</b>	<b>May 24, 2017</b>	<b>5</b>
3.1	Goals . . . . .	5
3.2	Progress . . . . .	5
3.3	EM Algorithm . . . . .	5
3.3.1	E step . . . . .	6
<b>4</b>	<b>May 27, 2017</b>	<b>6</b>
4.1	Progress . . . . .	6
4.2	EM Algorithm (cont.) . . . . .	6
4.2.1	M step . . . . .	6
<b>5</b>	<b>May 28, 2017</b>	<b>8</b>
5.1	Progress . . . . .	8
5.2	EM Algorithm (cont.) . . . . .	9
5.2.1	M step (cont.) . . . . .	9
5.3	Parametric Bootstrap . . . . .	9
<b>6</b>	<b>May 29, 2017</b>	<b>9</b>
6.1	Progress . . . . .	9
6.2	Goals . . . . .	9

<b>7</b>	<b>May 30, 2017</b>	<b>9</b>
7.1	Progress . . . . .	9
7.2	Goals . . . . .	10
7.3	Previous Model of EPSP Amplitudes - Binomial Distribution . . . . .	10
7.3.1	Statistical Model . . . . .	10
7.3.2	Method of Moments Estimator . . . . .	10
7.3.3	Maximum Likelihood Estimator . . . . .	11
7.3.4	Removing Zeros . . . . .	12
<b>8</b>	<b>May 31, 2017</b>	<b>12</b>
8.1	Progress . . . . .	12
8.2	Goals . . . . .	12
<b>9</b>	<b>June 1, 2017</b>	<b>12</b>
9.1	Progress . . . . .	12
9.2	Goals . . . . .	12
9.3	Meeting with Professor Barth and Professor Brasier . . . . .	13
<b>10</b>	<b>June 3, 2017</b>	<b>13</b>
10.1	Progress . . . . .	13
10.2	Goals . . . . .	13
10.3	Parsing IGOR files with R . . . . .	13
<b>11</b>	<b>June 5, 2017</b>	<b>13</b>
11.1	Progress . . . . .	13
<b>12</b>	<b>June 7, 2017</b>	<b>14</b>
12.1	Progress . . . . .	14
12.2	Goals . . . . .	14
12.3	Useful Properties of the Lognormal Distribution . . . . .	14
12.4	Parameter Estimates for the Lognormal Model . . . . .	14
<b>13</b>	<b>June 8, 2017</b>	<b>14</b>
13.1	Progress . . . . .	14
13.2	Goals . . . . .	15
13.3	Estimating Probabilities of Large Amplitudes . . . . .	15
<b>14</b>	<b>June 9, 2017</b>	<b>15</b>
14.1	Progress . . . . .	15
14.2	Goals . . . . .	15
<b>15</b>	<b>June 12, 2017</b>	<b>16</b>
15.1	Progress . . . . .	16
15.2	Goals . . . . .	16
<b>16</b>	<b>June 13, 2017</b>	<b>16</b>
16.1	Goals . . . . .	16

<b>17 June 14, 2017</b>	<b>16</b>
17.1 Progress . . . . .	16
17.2 Goals . . . . .	16
<b>18 June 15, 2017</b>	<b>16</b>
18.1 Progress . . . . .	16
18.2 Goals . . . . .	16
18.3 Meeting with Professor Barth and Professor Brasier . . . . .	17
<b>19 June 16, 2017</b>	<b>17</b>
19.1 Progress . . . . .	17
19.2 Goals . . . . .	17
<b>References</b>	<b>17</b>

# 1 May 20, 2017

## 1.1 Progress

- set up github for project
- plotted histograms of first trials
- plotted scatter plots of first trials (for stationarity)

## 1.2 Statistical Model

Recall that our model of the process is

$$X = \sum_{j=1}^N Z_j, \quad Y_j = \text{Bernoulli}(p_j), \quad (Z_j \mid Y_j = 1) \sim N(\mu_j, \sigma_j^2), \quad (Z_j \mid Y_j = 0) = 0$$

where  $N, \mu_j, \sigma_j^2, p_j$  are all unknown parameters of the model. The  $Z_j$  random variable models the response amplitude of a single contact of which there are  $N$ , and the  $Y_j$  random variable models the release success of a single contact. We assume that all of the  $Z_j$  and the  $Y_j$  are independent.

### 1.2.1 Justification

The additivity of the potential is justified by Petterson and Einevoll [6]. The use of a Gaussian distribution for individual contacts is justified by Magee and Cook [4].

# 2 May 22, 2017

## 2.1 Progress

- derived point estimates for release probability, mean response amplitude, and response amplitude variance under the constant parameter models

## 2.2 Point Estimation of Model Parameters

### 2.2.1 Estimations Under a Simplified Model

We have that the expectation of the model is

$$\mathbb{E}[X] = \sum_{j=1}^N \mathbb{E}[Z_j] = \sum_{j=1}^N \mu_j \cdot p_j$$

and that the variance is

$$\text{Var}[X] = \sum_{j=1}^N \text{Var}[Z_j] = \sum_{j=1}^N \sigma_j^2 \cdot p_j.$$

Also note that the failure rate of the model, i.e. the probability that all  $N$  contacts fail to release, can be approximated by

$$\mathbb{P}[X = 0] \approx \prod_{j=1}^N (1 - p_j)$$

by assuming that a contact produces a positive response everytime it succeeds in releasing a vesicle. Thus, with simplifying assumptions that all the  $\mu_j$ ,  $\sigma_j^2$ , and  $p_j$  are the same across the  $N$  points of contact, and by fixing a value of  $N$ , we may find a plugin estimator for  $\hat{p}$  and method of moments estimators for  $\hat{\mu}$  and  $\hat{\sigma}^2$  that depend on  $\hat{p}$ . If we let  $\bar{X}$  be the sample mean,  $S^2$  be the sample variance, and  $p_f$  be the sample failure rate, these estimates are given by

$$\hat{p} = 1 - \sqrt[N]{p_f}, \quad \hat{\mu} = \frac{\bar{X}}{N\hat{p}}, \quad \hat{\sigma}^2 = \frac{S^2}{N\hat{p}}.$$

### 2.2.2 EM Algorithm

Let us return to the general case. Suppose that we treat the response amplitudes of the individual contacts, the  $Z_j$  from  $1 \leq j \leq N$ , as latent variables. Then, the joint distribution of the latent and observed variables will be an exponential family, so EM algorithm should work very well.

## 3 May 24, 2017

### 3.1 Goals

- how was the data collected? can we really justify our model?
- workout details of EM and implement

### 3.2 Progress

- derive E step of the EM algorithm

### 3.3 EM Algorithm

We refer to lecture notes by Andrew Ng [5] for the EM reference. Suppose that we fix  $N$  and let  $(x_i)_{i=1}^n$  be our observations of our model, let  $y = (y_j)_{j=1}^N$  be the hidden observed values of the release success of each individual contact, and let  $(\mu, \sigma^2, p) = (\mu_j, \sigma_j^2, p_j)_{j=1}^N$ . Then, the log-likelihood of the parameters is given by

$$\ell(\mu, \sigma^2, p) = \sum_{i=1}^n \log p(x_i; \mu, \sigma^2, p) = \sum_{i=1}^n \log \sum_{y \in \{0,1\}^N} p(x_i | y; \mu, \sigma^2, p) \cdot p(y; \mu, \sigma^2, p).$$

In the above, the sum ranges over all possible assignments of the release success of the  $N$  individual contacts. Note that  $p(x_i | y; \mu, \sigma^2, p)$  is simply the pdf of the sum of Gaussian variables which also takes a Gaussian distribution, and  $p(y; \mu, \sigma^2, p)$  is simply given by the product of  $p_j$  or  $1 - p_j$  for each  $1 \leq j \leq N$ , as appropriate. Also note that this problem is very similar to a Gaussian mixture model, where we sample from  $2^N$  different Gaussian distributions where each Gaussian is a sum of

the original  $N$  Gaussians. The difference is that the Gaussians that we find may not be arbitrary, but in fact must be generated by  $N$  Gaussians in a specific way.

### 3.3.1 E step

In the E step, we find a probability distribution over the release success of the individual contacts, given the observed amplitude  $x_i$ . Denote  $\varphi(x; \mu, \sigma^2)$  as the pdf of a Gaussian with mean  $\mu$  and variance  $\sigma^2$  evaluated at the point  $x$ . Then,

$$Q_i(y) = p(y \mid x_i; \mu, \sigma^2, p) = \frac{p(x_i \mid y; \mu, \sigma^2, p) \cdot p(y; \mu, \sigma^2, p)}{\sum_{z \in \{0,1\}^N} p(x_i \mid z; \mu, \sigma^2, p) \cdot p(z; \mu, \sigma^2, p)}$$

where for all assignments of release successes  $z \in \{0, 1\}^N$ ,

$$p(x_i \mid z; \mu, \sigma^2, p) = \varphi\left(x_i; \sum_{j: z_j=1} \mu_j, \sum_{j: z_j=1} \sigma_j^2\right), \quad p(z; \mu, \sigma^2, p) = \prod_{j: z_j=1} p_j \prod_{j: z_j=0} (1 - p_j).$$

## 4 May 27, 2017

### 4.1 Progress

- derive M step of the EM algorithm for  $p_l$

### 4.2 EM Algorithm (cont.)

#### 4.2.1 M step

In the M step, we maximize the function

$$\ell(\theta) = \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \log \frac{p(x_i, z; \theta)}{Q_i^{(t)}(z)}$$

and set the value of  $\theta$  that maximizes the above to the new estimate for  $\theta$ , where  $\theta$  is our  $\mu, \sigma^2, p$ . For notational convenience, let  $\mu_z := \sum_{j: z_j=1} \mu_j$ ,  $\sigma_z^2 := \sum_{j: z_j=1} \sigma_j^2$ , and  $p_z := \prod_{j: z_j=1} p_j \prod_{j: z_j=0} (1 - p_j)$ .

- Maximization with respect to  $\mu_l$

We have that

$$\begin{aligned}
\frac{\partial}{\partial \mu_l} \ell(\theta) &= \frac{\partial}{\partial \mu_l} \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \log \frac{p(x_i, z; \theta)}{Q_i^{(t)}(z)} = \frac{\partial}{\partial \mu_l} \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \left( \log p(x_i, z; \theta) - \log Q_i^{(t)}(z) \right) \\
&= \frac{\partial}{\partial \mu_l} \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \left( \log \left( \frac{1}{\sqrt{2\pi\sigma_z^2}} e^{-\frac{(x_i - \mu_z)^2}{2\sigma_z^2}} \cdot p_z \right) - \log Q_i^{(t)}(z) \right) \\
&= \frac{\partial}{\partial \mu_l} \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \left( \log \frac{1}{\sqrt{2\pi\sigma_z^2}} - \frac{(x_i - \mu_z)^2}{2\sigma_z^2} + \log p_z - \log Q_i^{(t)}(z) \right) \\
&= -\frac{\partial}{\partial \mu_l} \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \frac{(x_i - \mu_z)^2}{2\sigma_z^2} = -\sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \frac{-2(x_i - \mu_z)}{2\sigma_z^2} \\
&= \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \frac{x_i - \mu_z}{\sigma_z^2}.
\end{aligned}$$

We cannot isolate  $\mu_l$ , so we need to solve with the other  $\mu_j$  as well as the  $\sigma_j^2$ .

- **Maximization with respect to  $\sigma_l^2$**

We have that

$$\begin{aligned}
\frac{\partial}{\partial \sigma_l^2} \ell(\theta) &= \frac{\partial}{\partial \sigma_l^2} \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \log \frac{p(x_i, z; \theta)}{Q_i^{(t)}(z)} \\
&= \frac{\partial}{\partial \sigma_l^2} \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \left( \log \frac{1}{\sqrt{2\pi\sigma_z^2}} - \frac{(x_i - \mu_z)^2}{2\sigma_z^2} + \log p_z - \log Q_i^{(t)}(z) \right) \\
&= \frac{\partial}{\partial \sigma_l^2} \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \left( \log \frac{1}{\sqrt{2\pi\sigma_z^2}} - \frac{(x_i - \mu_z)^2}{2\sigma_z^2} \right) \\
&= \frac{\partial}{\partial \sigma_l^2} \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \left( -\frac{\log(2\pi\sigma_z^2)}{2} - \frac{(x_i - \mu_z)^2}{2\sigma_z^2} \right) = \frac{1}{2} \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \left( \left( \frac{x_i - \mu_z}{\sigma_z^2} \right)^2 - \frac{1}{\sigma_z^2} \right).
\end{aligned}$$

Now we need to solve for all the  $\mu_j$  and  $\sigma_j$  when the above is set to 0. It is worth noting that the Z-score with respect to  $N(\mu_z, \sigma_z^2)$  appears in both of the derivatives for  $\mu_l$  and  $\sigma_l^2$ .

- **Maximization with respect to  $p_l$**

As before, we have that

$$\begin{aligned}
\frac{\partial}{\partial p_l} \ell(\theta) &= \frac{\partial}{\partial p_l} \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \log \frac{p(x_i, z; \theta)}{Q_i^{(t)}(z)} \\
&= \frac{\partial}{\partial p_l} \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \left( \log \frac{1}{\sqrt{2\pi\sigma_z^2}} - \frac{(x_i - \mu_z)^2}{2\sigma_z^2} + \log p_z - \log Q_i^{(t)}(z) \right) \\
&= \frac{\partial}{\partial p_l} \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \log p_z \\
&= \frac{\partial}{\partial p_l} \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \left( \sum_{j: z_j=1} \log p_j + \sum_{j: z_j=0} \log(1 - p_j) \right) \\
&= \frac{\partial}{\partial p_l} \sum_{i=1}^n \left( \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \log p_l + \sum_{\substack{z \in \{0,1\}^N \\ z_l=0}} Q_i^{(t)}(z) \log(1 - p_l) \right) \\
&= \sum_{i=1}^n \left( \frac{\sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z)}{p_l} - \frac{\sum_{\substack{z \in \{0,1\}^N \\ z_l=0}} Q_i^{(t)}(z)}{1 - p_l} \right).
\end{aligned}$$

Setting this to 0 gives us that

$$0 = \sum_{i=1}^n \left( \frac{\sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z)}{p_l} - \frac{\sum_{\substack{z \in \{0,1\}^N \\ z_l=0}} Q_i^{(t)}(z)}{1 - p_l} \right),$$

which solves to

$$p_l = \frac{\sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z)}{\sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) + \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=0}} Q_i^{(t)}(z)} = \boxed{\frac{\sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z)}{n}}.$$

## 5 May 28, 2017

### 5.1 Progress

- gave up on a closed form for the M step, decided on gradient descent for the EM algorithm



## 5.2 EM Algorithm (cont.)

### 5.2.1 M step (cont.)

Since we cannot solve for the optimal  $\mu_l$  and  $\sigma_l^2$  in closed form, we opt for numerically determining the  $\mu_l$  and  $\sigma_l^2$  via gradient descent [1]. Recall that we have already computed the gradient:

$$\begin{cases} \frac{\partial}{\partial \mu_l} \ell(\theta) = \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \frac{x_i - \mu_z}{\sigma_z^2} \\ \frac{\partial}{\partial \sigma_l^2} \ell(\theta) = \frac{1}{2} \sum_{i=1}^n \sum_{\substack{z \in \{0,1\}^N \\ z_l=1}} Q_i^{(t)}(z) \left( \left( \frac{x_i - \mu_z}{\sigma_z^2} \right)^2 - \frac{1}{\sigma_z^2} \right) \end{cases} .$$

Now, it is straightforward to implement the algorithm.

## 5.3 Parametric Bootstrap

Once we have an estimator for the parameters of the model, we may use parametric bootstrap to estimate the variance of the estimator, as explained in the classical text by Casella and Berger [2]. Using the parameters  $\hat{\theta}$  estimated by the EM algorithm, we may simulate  $n$  values  $X_1^*, \dots, X_n^* \sim f(x; \hat{\theta})$  and approximate the MLE using the EM algorithm  $B$  times, and use the sample variance of those  $B$  estimates of the MLE as the variance of the estimator.

## 6 May 29, 2017

### 6.1 Progress

- implement EM algorithm upto expectation step

### 6.2 Goals

- validate  $p_j$  estimates with the sample failure rate
- validate  $\mu_j$  estimates with the sample mean
- validate  $\sigma_j^2$  estimates with the sample variance
- after writing slower version of EM algorithm, write a faster version

## 7 May 30, 2017

### 7.1 Progress

- derive EM algorithm steps for “binomial distribution” model
- implement draft of EM algorithm for “binomial distribution” model
- test EM algorithm, realize that we need to handle the  $z = (0, 0, \dots, 0)$  case separately

## 7.2 Goals

- test EM algorithm with data simulated from the generative model
- justify linearity of postsynaptic integration of signals
- investigate emergent properties
- separate out nonzero responses from the zero responses
- show that the classical binomial model does not sufficiently explain the data (figures, graphics)
- investigate results of Turner and West and see if they apply
- also try model that keeps  $\sigma^2$  constant, and derive the EM algorithm for it
- study Kolmogorov-Smirnov test for testing the equality of the distributions

## 7.3 Previous Model of EPSP Amplitudes - Binomial Distribution

The previous models of postsynaptic potential amplitudes are usually binomial models [3]. We would like to first show that this model does not fit our data in a satisfactory way. This model assumes that the probabilities of all of the  $N$  contacts are the same and that the mean response is also all the same. Let us further suppose that the variances of the Gaussians at each contact are the same and see if we may reject this hypothesis. The insufficiency of this model is also suggested by Turner and West [7], and we wish to support this view.

### 7.3.1 Statistical Model

With  $N$  contacts, the model is given by

$$X = \sum_{j=1}^N Z_j, \quad Y_j = \text{Bernoulli}(p), \quad (Z_j \mid Y_j = 1) \sim N(\mu, \sigma^2), \quad (Z_j \mid Y_j = 0) = 0$$

with parameters  $\mu, \sigma^2, p$ .

### 7.3.2 Method of Moments Estimator

Recall that

$$\mathbb{E}[X] = \sum_{j=1}^N \mathbb{E}[Z_j] = Np\mu, \quad \text{Var}[X] = \sum_{j=1}^N \text{Var}[Z_j] = Np\sigma^2.$$

Then if we determine  $p$  via the failure rate  $p_f$ , we find estimates

$$\hat{p} = 1 - \sqrt[p_f]{p_f}, \quad \hat{\mu} = \frac{\bar{X}}{N\hat{p}}, \quad \hat{\sigma}^2 = \frac{S^2}{N\hat{p}},$$

as before.

### 7.3.3 Maximum Likelihood Estimator

The log likelihood of the above model as a function of the parameters is given by

$$\ell(\mu, \sigma^2, p) = \sum_{i=1}^n p(x_i; \mu, \sigma^2, p) = \sum_{i=1}^n \log \sum_{z \in \{0,1\}^N} p(x_i, z; \mu, \sigma^2, p).$$

As above, we must employ the EM algorithm to estimate the MLE. Under these assumptions, the E step simplifies to

$$Q_i(y) = p(y \mid x_i; \mu, \sigma^2, p) = \frac{p(x_i \mid y; \mu, \sigma^2, p) \cdot p(y; \mu, \sigma^2, p)}{\sum_{z \in \{0,1\}^N} p(x_i \mid z; \mu, \sigma^2, p) \cdot p(z; \mu, \sigma^2, p)}$$

with

$$p(x_i \mid z; \mu, \sigma^2, p) = \varphi(x_i; N_z \mu, N_z \sigma^2), \quad p(z; \mu, \sigma^2, p) = p^{N_z} (1-p)^{N-N_z} = p^{N_z} - p^N.$$

where  $N_z$  denotes the number of 1s in  $z \in \{0,1\}^N$ . In the M step, we have now have closed form solutions to all of the parameters. We have that

$$\frac{\partial}{\partial \mu} \ell(\theta) = \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \frac{-(x_i - N_z \mu)(-2N_z)}{2N_z \sigma^2} = \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \frac{x_i - N_z \mu}{\sigma^2} = 0$$

which gives a solution of

$$\mu = \frac{\sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i(z) x_i}{\sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i(z) N_z} = \boxed{\frac{\sum_{i=1}^n x_i}{\sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) N_z}}$$

for  $\mu$ ,

$$\begin{aligned} \frac{\partial}{\partial \sigma^2} \ell(\theta) &= \frac{\partial}{\partial \sigma^2} \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \left( -\frac{\log(2\pi N_z \sigma^2)}{2} - \frac{(x_i - N_z \mu)^2}{2N_z \sigma^2} \right) \\ &= \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \left( -\frac{1}{2\sigma^2} + \frac{1}{2} \left( \frac{x_i - N_z \mu}{N_z \sigma^2} \right)^2 \right) = 0 \end{aligned}$$

which gives a solution of

$$\sigma^2 = \frac{\sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \left( \frac{x_i - N_z \mu}{N_z} \right)^2}{\sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z)} = \boxed{\frac{\sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \left( \frac{x_i - N_z \mu}{N_z} \right)^2}{n}}$$

for  $\sigma^2$ , and

$$\begin{aligned} \frac{\partial}{\partial p} \ell(\theta) &= \frac{\partial}{\partial p} \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) (N_z \log p + (N - N_z) \log(1-p)) \\ &= \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \left( \frac{N_z}{p} - \frac{N - N_z}{1-p} \right) = 0 \end{aligned}$$

which gives a solution of

$$p = \frac{\sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) N_z}{\sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) N} = \boxed{\frac{\sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) N_z}{Nn}}$$

for  $p$ .

### 7.3.4 Removing Zeros

The above algorithms for EM in fact do not work, since they assign infinite probabilities to the 0s when they are assigned to the all failure case. Thus, we must implement the algorithm on the nonzero values. In this case, the  $z$  are drawn from  $S := \{0,1\}^N \setminus \{(0, \dots, 0)\}$  uniformly, and each  $z \in S$  is drawn with probability

$$\frac{1}{1 - \prod_{j=1}^N (1 - p_j)} \prod_{j: z_j=1} p_j \prod_{j: z_j=0} (1 - p_j).$$

In the case of the binomial distribution model, the partial with respect to  $p$  then turns out to be

$$\begin{aligned} \frac{\partial}{\partial p_l} \ell(\theta) &= \frac{\partial}{\partial p_l} \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \log p_z \\ &= \frac{\partial}{\partial p_l} \sum_{i=1}^n \sum_{z \in \{0,1\}^N} Q_i^{(t)}(z) \left( \sum_{j: z_j=1} \log p_j + \sum_{j: z_j=0} \log(1 - p_j) - \log \left( 1 - \prod_{j=1}^N (1 - p_j) \right) \right) \end{aligned}$$

## 8 May 31, 2017

### 8.1 Progress

- start making visualizations for MOMÉ estimators

### 8.2 Goals

- get more data!

## 9 June 1, 2017

### 9.1 Progress

- clear up problems on current approaches (lognormal distribution, lack of data)

### 9.2 Goals

- explore other distributions for the postsynaptic response at each synaptic contact
- use lognormal distribution at each contact instead of normal distribution

### 9.3 Meeting with Professor Barth and Professor Brasier

The immediate goal at the moment is to do the first figures of the paper, e.g. the statement of the problem, why this small data set is enough to suspect that a simple binomial model will not suffice to explain the data.

- failure rates of the trials, analysis with simple binomial model
- traces of the data
- aggregate results

## 10 June 3, 2017

### 10.1 Progress

- found IGOR file parser <https://www.rdocumentation.org/packages/IgorR>, plotted wave files with it

### 10.2 Goals

- view IGOR files on R
- ggplot histograms of the data, whole and first col
- failure rates of the columns (trials)
- plot failure rates of the columns as a function of the trials
- estimate probabilities of the highest amplitudes using binomial model, Markov, and Chebyshev
- awesome legends in the plots please (toggle with true/false options)
- plotting aggregate data
  - other ways to represent the whole data set

### 10.3 Parsing IGOR files with R

Use the `IgorR` package for reading the `pxp` files. Read the file with `read.pxp` and find the appropriate wave data. Then, convert to a `ts` object via `WaveToTimeSeries`. Finally, to plot with `ggplot2`, turn it into a data frame via `data.frame(Y = as.matrix(sweep1), t = time(sweep1))`.

## 11 June 5, 2017

### 11.1 Progress

- ggplot histograms, first col and all
- plot failure rates, individually and aggregate

## 12 June 7, 2017

### 12.1 Progress

- updated parameter estimates for the compound binomial-lognormal model

### 12.2 Goals

- check model with other predictions, such as mean, variance, and failure rate

### 12.3 Useful Properties of the Lognormal Distribution

The cdf and quantile functions are given by

$$\begin{cases} F(x) = \Phi \left[ \frac{\log(x) - \mu}{\sigma} \right] \\ F^{-1}(p) = \exp \left[ \mu + \sigma \Phi^{-1}(p) \right] \end{cases} \quad p \in (0, 1)$$

The  $t$ th moment is given by

$$\mathbb{E} [X^k] = \exp \left[ \mu k + \frac{\sigma^2 k^2}{2} \right].$$

In particular, the mean and variance are given by

$$\begin{cases} \mathbb{E}[X] = \exp \left[ \mu + \frac{\sigma^2}{2} \right] \\ \text{Var}[X] = \exp [2(\mu + \sigma^2)] - \exp [2\mu + \sigma^2] \end{cases}.$$

### 12.4 Parameter Estimates for the Lognormal Model

Under our new model that uses the lognormal distribution, the  $p$  parameter estimate can stay, but we need to fix the estimates for  $\mu$  and  $\sigma$ . Given  $N$  and  $p$ , we have that

$$\begin{cases} \mathbb{E}[X] = \sum_{j=1}^N \mathbb{E}[Z_j] = N \exp \left[ \mu + \frac{\sigma^2}{2} \right] p \\ \text{Var}[X] = \sum_{j=1}^N \text{Var}[Z_j] = N \left( \exp [2(\mu + \sigma^2)] - \exp [2\mu + \sigma^2] \right) p \end{cases}.$$

Now estimate  $\mathbb{E}[X]$  by the sample mean  $\bar{X}$  and  $\text{Var}[X]$  by the sample variance  $S^2$ . Then, after massaging, we get that the method of moments parameter estimates are

$$\begin{pmatrix} \mu \\ \sigma \end{pmatrix} = \begin{pmatrix} -1 & 1 \\ 2 & -1 \end{pmatrix} \begin{pmatrix} \frac{1}{2} \log \left( \frac{S^2}{Np} + \left( \frac{\bar{X}}{Np} \right)^2 \right) \\ 2 \log \left( \frac{\bar{X}}{Np} \right) \end{pmatrix}.$$

## 13 June 8, 2017

### 13.1 Progress

- constructed bootstrap pivot confidence intervals for binomial parameters estimates
- plotted bootstrap results for each test
- plotted bootstrap results across tests

## 13.2 Goals

- construct pivot confidence intervals for the binomial parameters
- plot bootstrap results for each parameter in a separate plot
- plot probability of large events (Markov, Chebyshev)
- add legends to testing plots (low priority)

## 13.3 Estimating Probabilities of Large Amplitudes

We use the Markov and Chebyshev inequalities to bound the tail probabilities. Markov's inequality gives us that

$$\mathbb{P}[X \geq t] \leq \frac{\mathbb{E}[X]}{t}$$

and Chebyshev's inequality gives us that

$$\mathbb{P}[|X - \mathbb{E}[X]| \geq t] \leq \frac{\text{Var}[X]}{t^2}.$$

Note that both

$$\mathbb{E}[X] = N \exp \left[ \mu + \frac{\sigma^2}{2} \right] p$$

and

$$\text{Var}[X] = N \left( \exp [2(\mu + \sigma^2)] - \exp [2\mu + \sigma^2] \right) p = N p e^{2\mu + \sigma^2} \left( e^{\sigma^2} - 1 \right)$$

are increasing functions of the parameters. Thus, we may find looser estimates that have more confidence by choosing the upper bound of the confidence intervals for the parameters.

## 14 June 9, 2017

### 14.1 Progress

- plotted upper bounds for simulations
- plotted cumulative upper bounds
- plotted upper bounds for the max amplitude (Markov, Chebyshev)
- plotted upper bounds for the max amplitude (Monte Carlo)

### 14.2 Goals

- plot for each type of theta (actual, estimated, upper bound)
- apply tests to actual data now
- approximation of probability of large amplitudes by simulation

## **15 June 12, 2017**

### **15.1 Progress**

- plotted Monte Carlo results for just the first column

### **15.2 Goals**

- plot Monte Carlo results for just the first column

## **16 June 13, 2017**

### **16.1 Goals**

- simulated histogram vs the original histogram
- finish figure 1

## **17 June 14, 2017**

### **17.1 Progress**

- made slides and photoshop pictures

### **17.2 Goals**

- are high amplitudes correlated with each other for a single sweep?
- work on slides
- draw figures with photoshop

## **18 June 15, 2017**

### **18.1 Progress**

- organize figure 1 ideas

### **18.2 Goals**

- figure 1



### 18.3 Meeting with Professor Barth and Professor Brasier

First figure should consist of something like:

- a. reconstructions
- b. multiple sweeps (maybe 10?) on top of each other
- c. heat map of amplitudes
- d. failure rates
- e. average failure rates

## 19 June 16, 2017

### 19.1 Progress

- plotted average failure rate plot (with and without ci)
- plotted heatmaps

### 19.2 Goals

- 

## References

- [1] Christopher M Bishop. Pattern recognition. *Machine Learning*, 128:1–58, 2006.
- [2] George Casella and Roger L Berger. *Statistical inference*, volume 2. Duxbury Pacific Grove, CA, 2002.
- [3] U Kuhnt, G Hess, and LL Voronin. Statistical analysis of long-term potentiation of large excitatory postsynaptic potentials recorded in guinea pig hippocampal slices: binomial model. *Experimental brain research*, 89(2):265–274, 1992.
- [4] Jeffrey C Magee and Erik P Cook. Somatic epsp amplitude is independent of synapse location in hippocampal pyramidal neurons. *Nature neuroscience*, 3(9):895–903, 2000.
- [5] Andrew Ng. The em algorithm. University Lecture, 2016.
- [6] Klas H Pettersen and Gaute T Einevoll. Amplitude variability and extracellular low-pass filtering of neuronal spikes. *Biophysical journal*, 94(3):784–802, 2008.
- [7] Dennis A Turner and Mike West. Bayesian analysis of mixtures applied to post-synaptic potential fluctuations. *Journal of neuroscience methods*, 47(1):1–21, 1993.