

---

## Problem Set 7

This problem set is due in class on **Tuesday, April 23, 1:30pm**. Please show your work and turn in all Matlab code and plots.

If you have questions, please post them on the Discussion Board on the Blackboard course website, rather than emailing the course staff. This will allow other students with the same question to see the response and any ensuing discussion.

The Matlab code referred to in the following problems is posted on the Blackboard course website under “Course Content → Data sets → ps7\_code.zip”. The dataset is posted under “Course Content → Data sets → ps7\_data.mat”. When you load the .mat file, you will find the following variables:

**Spikes:** a  $31 \times 552$  matrix of spike snippets<sup>1</sup>, where `Spikes(:,n)` is the  $n$ th snippet ( $n = 1, \dots, 552$ ). Values are in  $\mu\text{V}$ .

**InitParams:** a structure containing initialization parameters for a Gaussian mixture model, with the following fields

- **mu** is a  $2 \times 8$  matrix, where the  $k$ th column is the initialization of the  $k$ th cluster center,  $\mu_k$  ( $k = 1, \dots, K$ ).
- **Sigma** is a  $2 \times 2$  covariance matrix. Assume that all cluster covariances  $\Sigma_k$  are initialized to the same covariance matrix.

In Problem Sets 5 and 6, we clustered the 31-dimensional spike waveforms directly. Here, we will first use PCA to project each 31-dimensional waveform down to a two-dimensional feature space. Then, we will perform model selection in the feature space.

### 1. Automatic feature selection using PCA

Treat each snippet as a point  $\mathbf{x}_n \in \mathbf{R}^D$  ( $n = 1, \dots, N$ ), where  $D = 31$  is the number of samples in each snippet and  $N = 552$  is the number of detected spikes.

- (a) **(2 points)** Plot the  $N$  raw spike snippets in a “voltage vs. time” plot, as in Figure 7(a) of the Lewicki paper.
- (b) **(8 points)** Apply PCA to all  $N$  spike snippets. Plot the eigenvector waveforms corresponding to the three largest eigenvalues, as in Figure 7(b) of the Lewicki paper. Label the waveforms by color (1st: red, 2nd: green, 3rd: blue).

---

<sup>1</sup>The neural data have been generously provided by the laboratory of Prof. Krishna Shenoy at Stanford University. The data are to be used exclusively for educational purposes in this course.

- (c) **(2 points)** Plot the square-rooted eigenvalue spectrum, as in Figure 7(c) of the Lewicki paper. Look for an elbow in the eigenvalue spectrum. How many dominant eigenvalues are there?
- (d) **(5 points)** Create a scatter plot of the PC1 score versus the PC2 score, where each point corresponds to a spike snippet, as in Figure 7(d) of the Lewicki paper. How many distinct clusters do you see in the plot?

## 2. Model selection for number of clusters

In Problem 1(d), each spike snippet is represented as a point in a two-dimensional feature space found by PCA. We will now perform clustering using a Gaussian mixture model (GMM) in this two-dimensional space. For fitting a GMM using EM, you can either use your own code from Problem Set 6 or our code (`func_GMM.m`).

As discussed in class, the PC directions are only unique up to a sign difference. Before doing this problem, please check that your answer to Problem 1(d) has the same sign orientation as in `ps7_1d.jpg`, since we've chosen the EM initialization parameters to be sensible for this sign orientation.

- (a) **(20 points)** Compute the cross-validated likelihoods for a GMM applied to the two-dimensional PCA projections for  $K = 1, \dots, 8$ . For each value of  $K$ , perform four-fold cross-validation by dividing the dataset into four equally-sized partitions (i.e., fold 1 has  $n = 1, \dots, 138$ ; fold 2 has  $n = 139, \dots, 276$ ; ...).

Initialize EM using the parameters in `InitParams`. To initialize the  $\mu_k$ , use the first  $K$  columns of `InitParams.mu`. To initialize the  $\Sigma_k$ , use the same `InitParams.Sigma` for each cluster. To initialize the  $\pi_k$ , use  $1/K$  for each cluster.

Plot the cross-validated likelihoods versus  $K$ . What is the optimal value of  $K$ ? (Hint: See Section 4.6 in the Lewicki paper for an explanation of why the optimal value of  $K$  may be larger than you expected.)

- (b) **(8 points)** For each value of  $K = 1, \dots, 8$ , create a separate plot with:
  - the data scatter of PC1 score versus PC2 score, as in Problem 1(d) (this will be the same for each value of  $K$ )
  - a one-standard-deviation ellipse for each cluster ( $k = 1, \dots, K$ ) based on  $\Sigma_k$  centered at  $\mu_k$ . Use the parameters obtained from the first cross-validation fold. To plot ellipses, use `func_plotEllipse.m`.
- (c) **(5 points)** For  $K = 3$ , plot the canonical spike waveform corresponding to each cluster center in a “voltage versus time” plot. This will involve projecting the two-dimensional  $\mu_k$  out into the 31-dimensional space. Use the  $\mu_k$  from the first cross-validation fold.