

1. Given a d -dimensional mean vector \mathbf{v} and $d \times d$ covariance matrix \mathbf{C} (symmetric, pos def), generate N random sample points distributed according to a Gaussian with mean \mathbf{v} and covariance \mathbf{C} . Recall that we discussed in class how to do this by first using `randn` in matlab to generate variables \mathbf{x} with zero mean and covariance \mathbf{I} , and then doing a linear transformation $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ to transform them into new variables \mathbf{y} with the desired mean and covariance.

2. For a set of 2D points generated as above, make a 2D plot of them in a graphics display, and then draw overlaid on that an elliptical contour associated with a Gaussian of mean \mathbf{v} and covariance \mathbf{C} to verify that your transformed points form a cluster with the correct center location and elliptical spread. We discussed doing this by first generating points on a 2D circle with radius 1.5 and then doing a linear transformation $\mathbf{A}\mathbf{x} + \mathbf{b}$ to transform those points into new points lying on an elliptical boundary.

Note: as a sanity check, try all three cases we talked about in class. That is, try to generate and draw data/contours from: 1) circular Gaussian, 2) elliptical Gaussian with major axis lying along either the x or y coordinate axis, and 3) rotated elliptical Gaussian (major axis is oriented diagonally).

3. Given a set of N sample points from a Gaussian distribution with d -dimensional mean vector \mathbf{v} and $d \times d$ covariance matrix \mathbf{C} , compute the maximum likelihood estimates for \mathbf{v} and \mathbf{C} using the formulas from our class notes. If you generate the sample points using the routine in part 1 and some given mean \mathbf{v}_0 and covariance \mathbf{C}_0 , and you then estimate the sample mean \mathbf{v} and \mathbf{C} using the maximum likelihood routine, you might expect \mathbf{v} to be close to \mathbf{v}_0 (distance-wise), and \mathbf{C} to be similar to \mathbf{C}_0 (similar axes and eigenvalues).

4. Given the specification of a mixture of Gaussians distribution with K Gaussian components, generate N sample points from that distribution. The values that are given to you are the mixing weights p_1, p_2, \dots, p_K , the mean vector of each component $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_K$, and the covariance matrix of each component $\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_K$. Recall from our class the generative process for a mixture of Gaussians, which can easily be turned into an algorithm: for $n=1$ to N , generate sample point \mathbf{x}_n by first choosing which of the K components it should come from (with the help of a uniform random number generator), and then using the routine from part 1 to generate a sample point from that component.

5. Given a set of N sample points from a mixture of Gaussians distribution with a known number of components K , use the EM algorithm to estimate the values of the parameters of the mixture of Gaussians (that is, the K mixing weights, the K mean vectors, and the K covariance matrices). If you run this for 2-dimensional data points, you should be able to visualize the results of this step by plotting the data points overlaid with covariance ellipses for each Gaussian component. If you generate a sample set from part 4, and then estimate the parameters using EM, you should get back parameters that are "close" to the true parameters used to generate the data.

6. Try your EM estimation routine on some real data. I will post some ideas on the web site, but you are welcome (indeed, encouraged) to try it on data that has some meaning to you.