

CS 274A Homework 5

Zachary DeStefano, 15247592

Due Date: Wednesday March 5th

K-means plots

Plots of Clusters

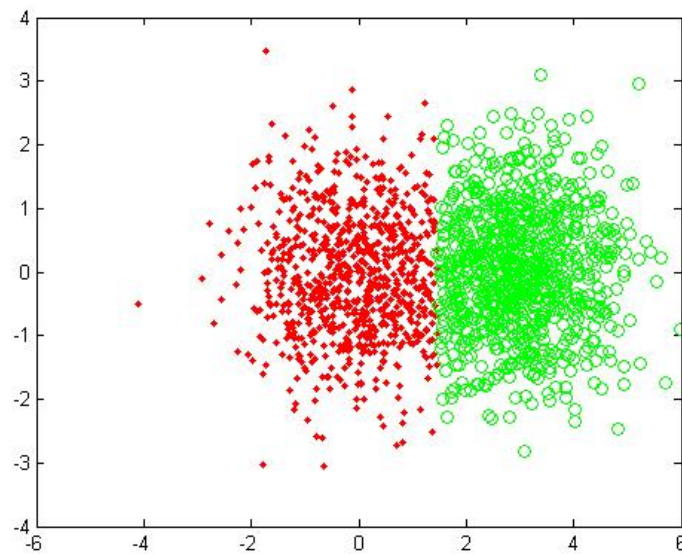


Figure 1: The k-Means plot for dataset1, $K=2$

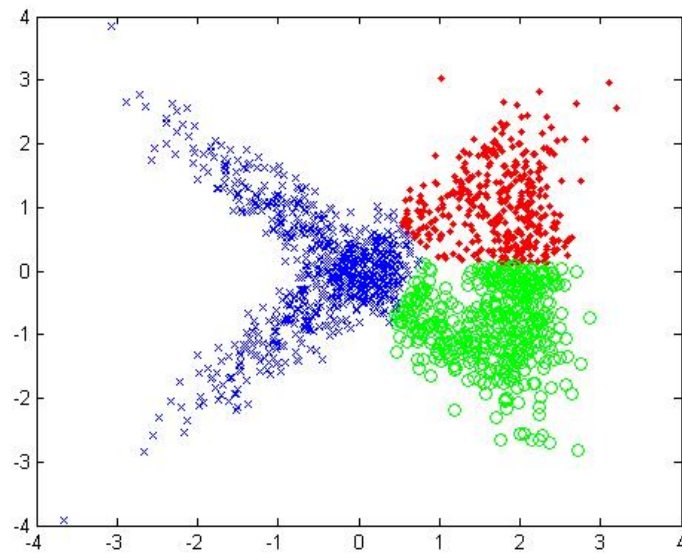
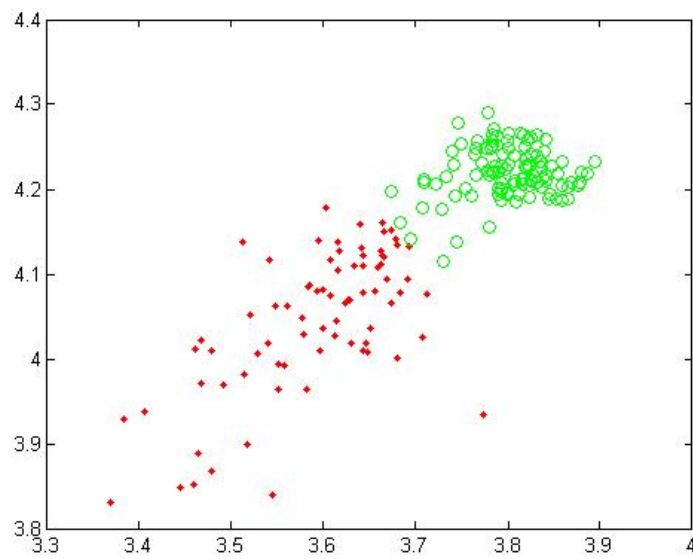
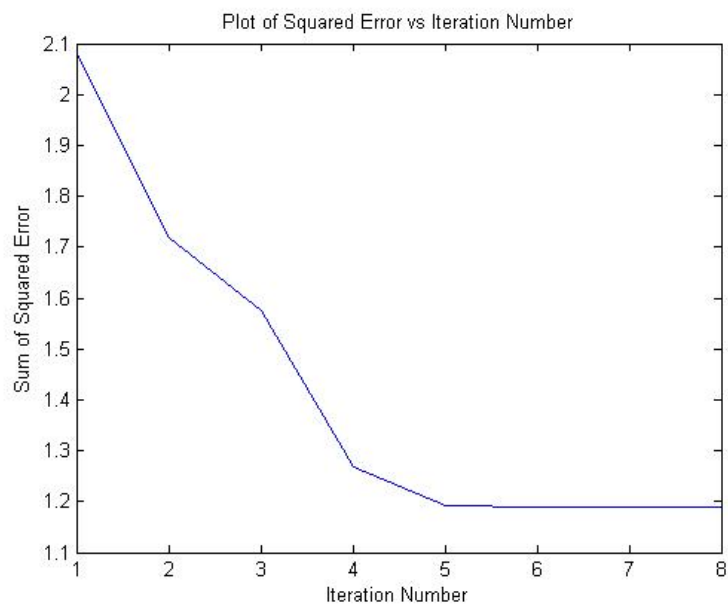
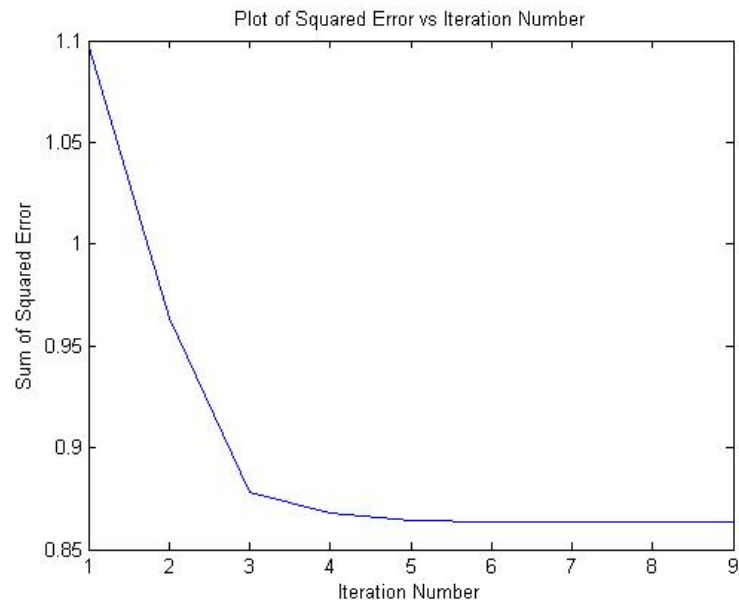
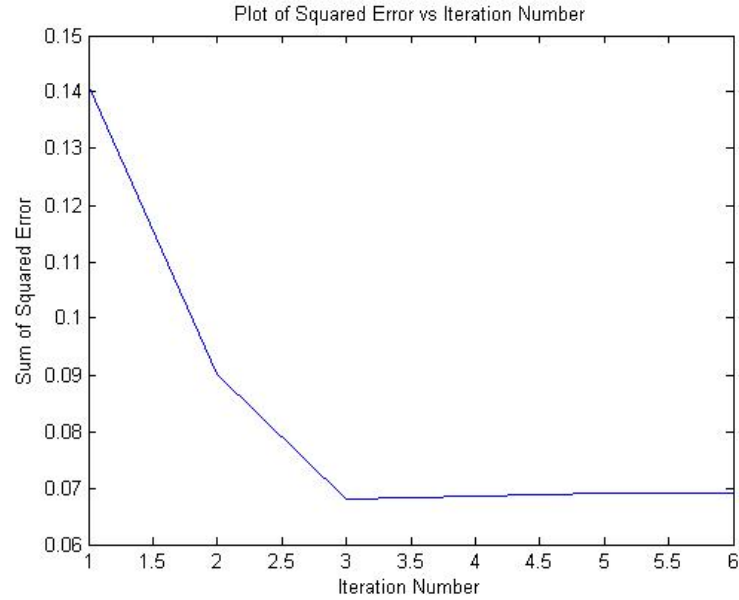


Figure 2: The k-Means plot for dataset2, $K=3$

Figure 3: The k-Means plot for dataset3, $K=2$

Plots of Sum of Squared Errors versus Iteration

Figure 4: The Sum of Squared Error plot for dataset1, $K=2$

Figure 5: The Sum of Squared Error plot for dataset2, $K=3$ Figure 6: The Sum of Squared Error plot for dataset3, $K=2$

EM plots

Plots of clusters using EM algorithm

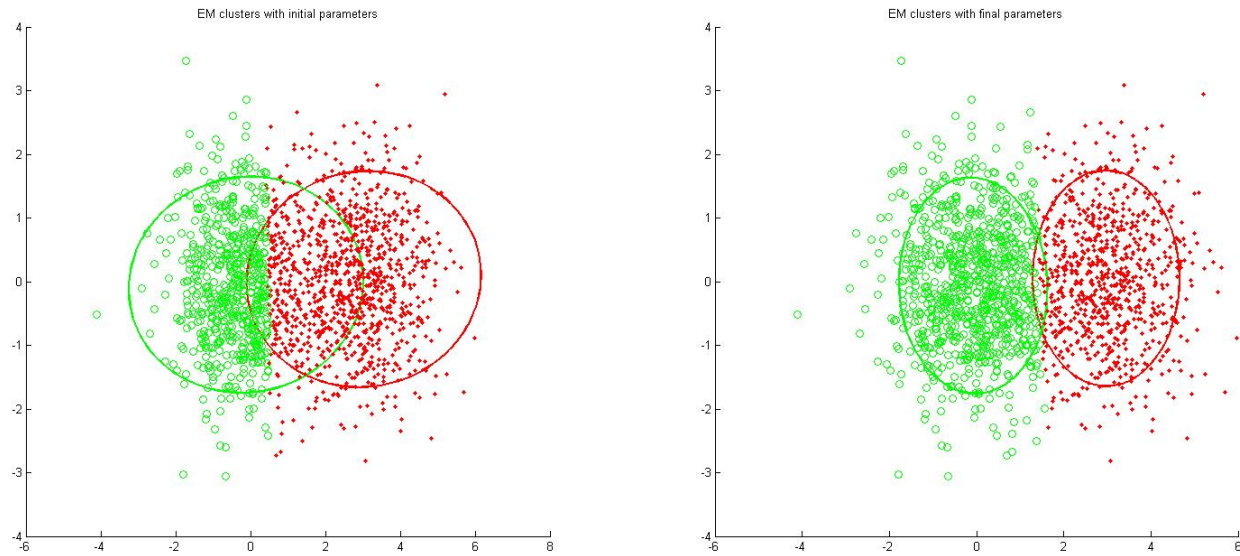


Figure 7: The EM cluster plots for dataset1, $K=2$, using initialization method 3

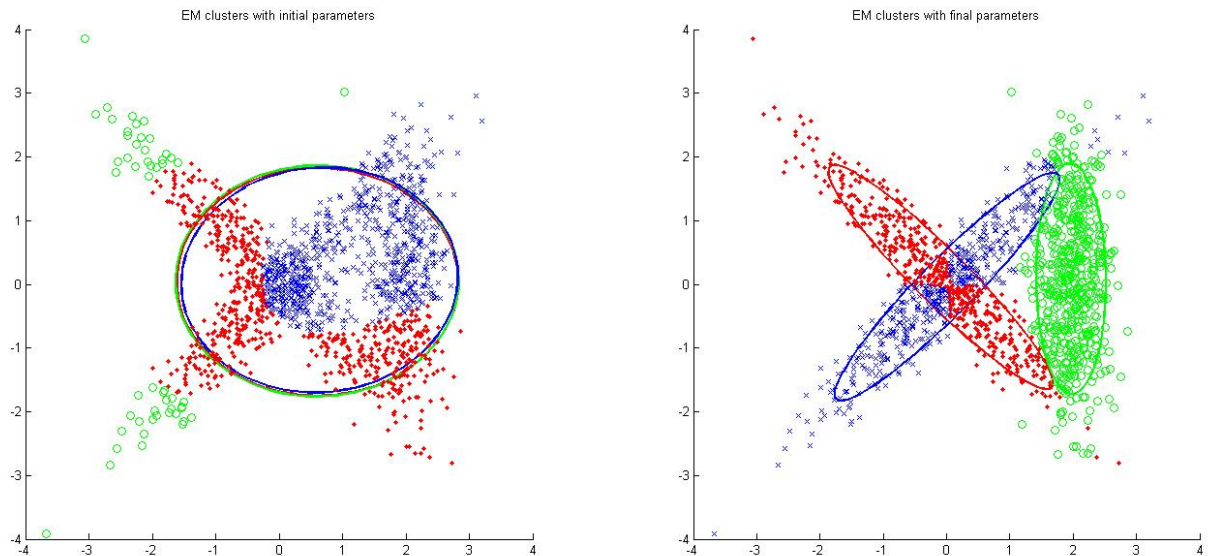


Figure 8: The EM cluster plots for dataset2, $K=3$, using initialization method 1

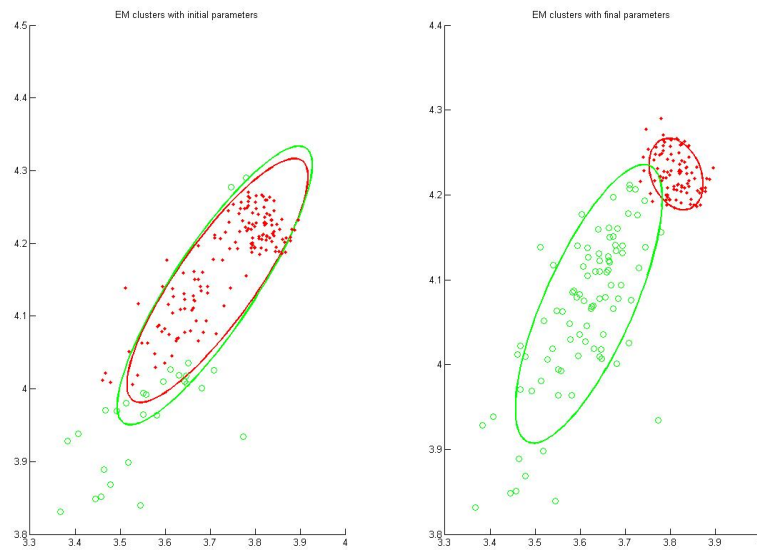


Figure 9: The EM cluster plots for dataset3, $K=2$, using initialization method 1

Plots of likelihood using EM algorithm

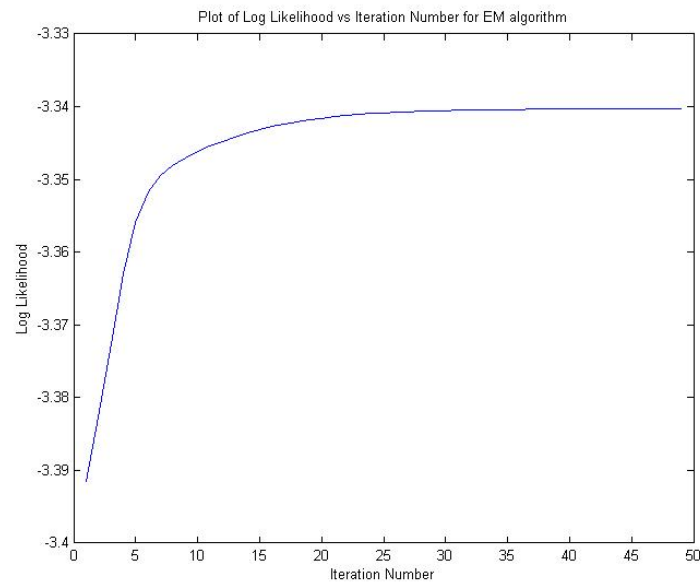


Figure 10: The likelihood plot for dataset1, $K=2$, using initialization method 3

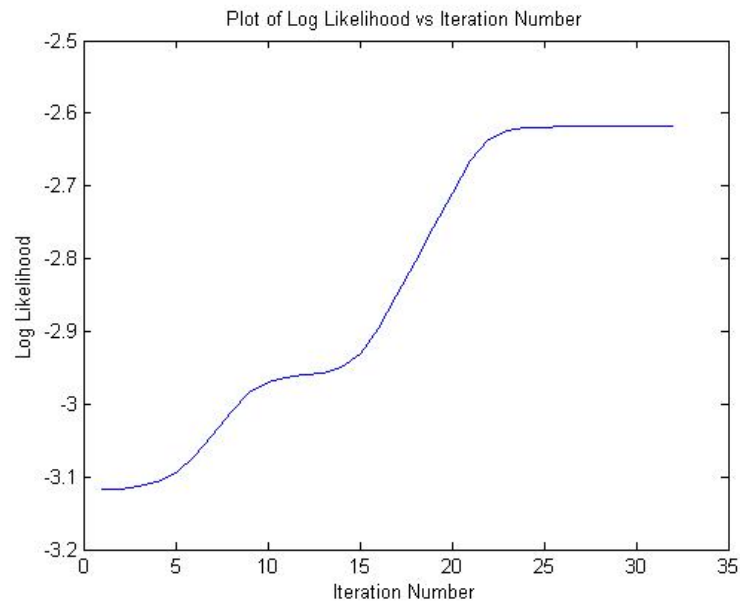


Figure 11: The likelihood plot for dataset2, $K=3$, using initialization method 1

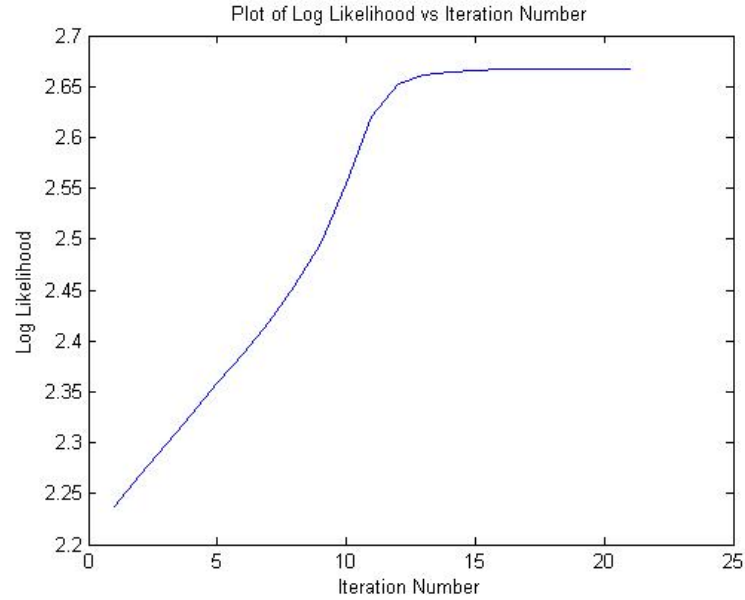


Figure 12: The likelihood plot for dataset3, $K=2$, using initialization method 1

Comments on Results

For dataset1, there was not much of a difference between the K-means result and the EM algorithm result. This was likely due to the fact that the two Gaussians did not overlap so K-means put the points in the correct Gaussian cluster more easily. With dataset2, there was a substantial difference. The EM algorithm seemed to capture the Gaussians much better than K-means. This is likely due to the fact that two of the Gaussians had an overlapping mean. Since K-means just clusters them into groups, it would not detect the differing Gaussians, whereas since EM is specifically trying to fit to Gaussians, it would detect it. For dataset3, the algorithm performed slightly better. ****INSERT COMMENT ABOUT RESULT WHEN COMPARING WITH LABELSET3****

In the end, because this data was generated from Gaussian densities and the EM algorithm fits the data to Gaussians, the EM algorithm was ideally suited to find the clusters in our case. However, if our need was to put data into two groups, as with dataset3, without much knowledge of the underlying model, then k-means is the ideal solution. ****INSERT COMMENT AFTER RUNNING COMPARISON WITH DATASET3****