Alex Ho
ID: 1001294712
CSE 4334-001
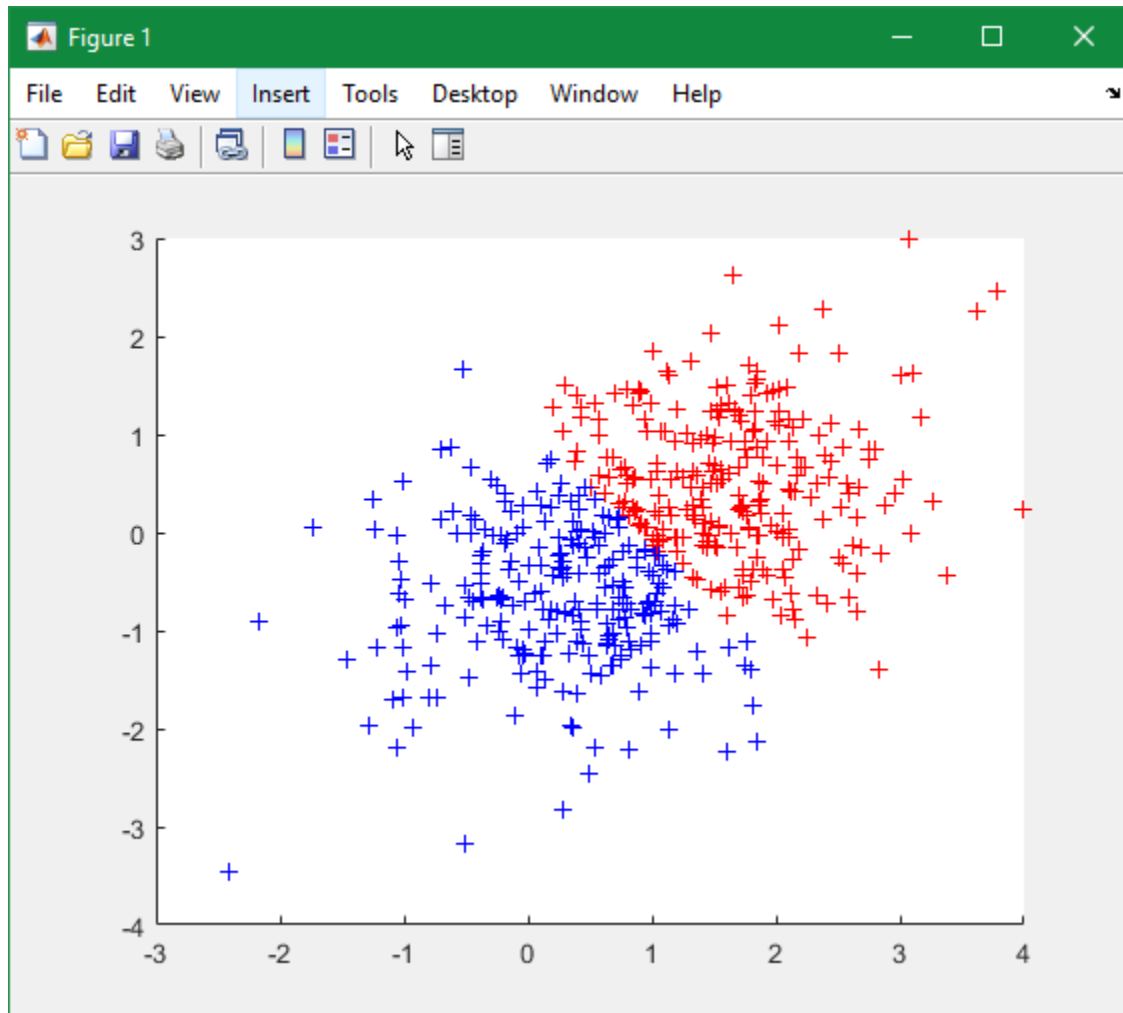
Assignment 1 Report

**Problem 1:**

Part 1:

Mykmeans.m was relatively simple to implement. It completes k means and returns an array of the data point's clusters. However, since it's only returning the centers they belong to, I wasn't sure how to plot the new centers outside of the method. I do have debug code within the method that will plot it, but I left it commented.

The biggest problem I had when solving the problem was figuring out how to converge the centers, but I was able to do it by getting the average of all the data points and setting the new center to it. My function also has two built in functions, one that finds the L2 from a point to a given center vector, and one that takes two vectors and finds the distances for each corresponding row. The first one is used when iterating through each data point to find its distance to every center. The second one was used to simply find the distance from all old centers to all new centers.

Part 2:

Applying the code to the data given with the specified parameters in part 2, the run in Figure 1 below took 11 iterations. The centers found are at (1.6737, 0.5212) and (0.2181, -0.6804). I did not plot the centers in the figure, as it would require plotting in the method itself.
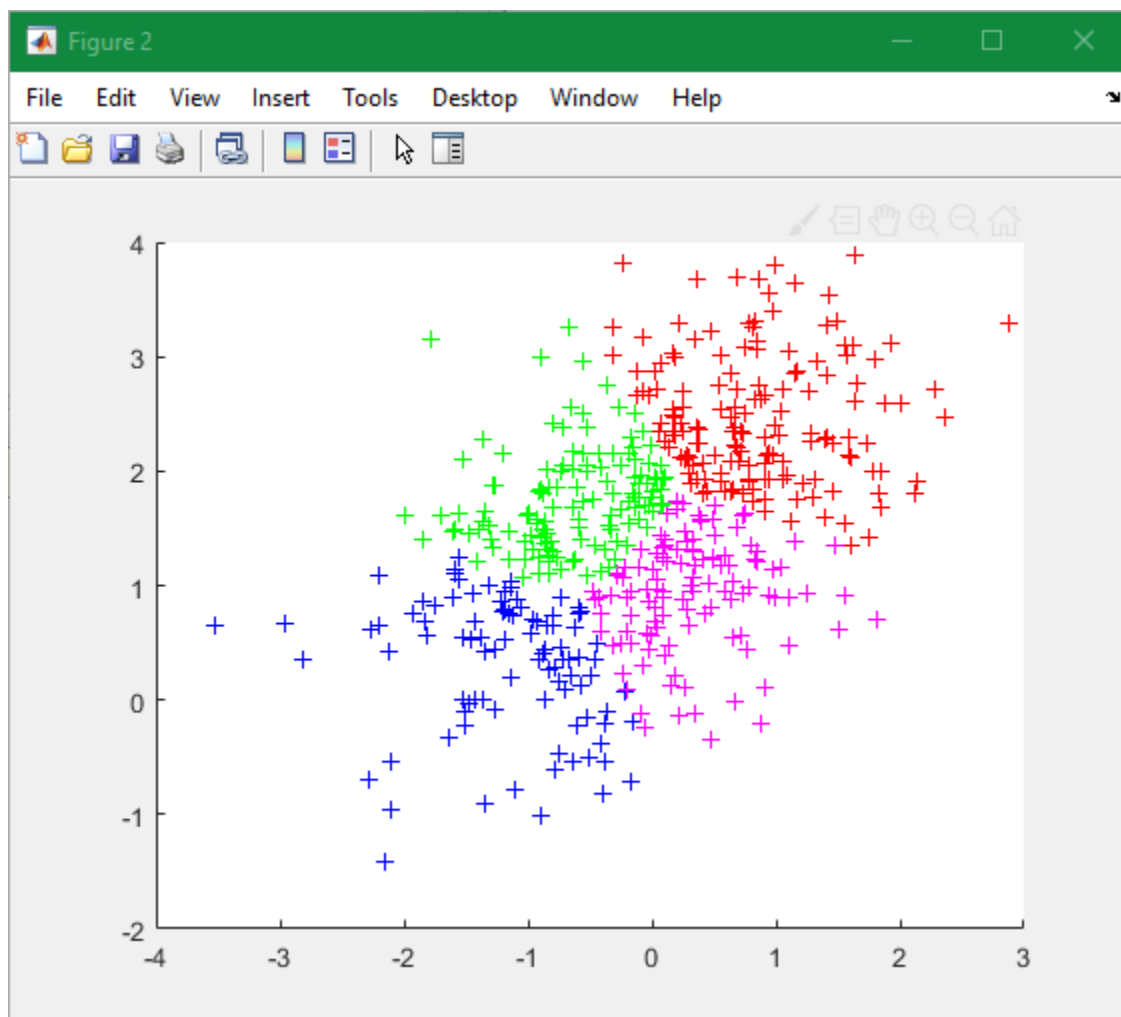
**Figure 1**



Part 3:

Doing the same for part 3 as I did in part 2, it took 9 iterations. The centers were at (.8642, 2.4646) , (-1.1760, 0.3076), (0.3224, 0.9326) and (-0.6405, 1.7524). The result is in Figure 2

**Figure 2**

**Problem 2:**

Problem 2 was much harder. I was confused on how to implement the kernel density estimation formula, but I was able to figure it out in the end. However, there are a couple of potential issues with the code:

1: I wasn't entirely sure what kernel function I should use. I decided to go with the Parzen Window function in the slides, but I also included a Gaussian kernel function commented out within the method.

2: I didn't understand what data Part 3 was asking for. It listed two datasets, but the first dataset was the same as Part 1. In the end, I decided to omit the duplicate data set and only use the new new mu and sigma (0 and 0.2).

3: When I compared mykde to the built in kde function, I noticed that their line reached all the way to the bottom. This goes past the given domain, but I wasn't sure how to implement it or if I was supposed to. Therefore, my kde line stops at the end of the domain.

4: I had trouble understanding how to implement a the multivariate kernel density estimation for a matrix. I ended up using the built in mvksdensity function to compute the probability.

Beyond these problems, the function should work as intended. The figures provided below corrispond these parts:

**Figure 1-4** = kde of Part 2 at bandwidths h = {.1, 1, 5, 10}. Data is N = 1000 gaussian random data with $\mu = 5$ and $\sigma 1 = 1$

**Figure 5-8** = kde of Part 3 at bandwidths h = {.1, 1, 5, 10}. Data is N = 1000 gaussian random data with $\mu = 0$ and $\sigma 1 = 0.2$

**Figure 9-12** = kde of Part 4, $N_1$ at bandwidths h = {.1, 1, 5, 10}. Data is N1 = 500 gaussian random data with $\mu = [1, 0]$ and $\Sigma = [0.9, 0.4 ; 0.4 , 0.9]$.

**Figure 13-16** = kde of Part 4, $N_2$ at bandwidths h = {.1, 1, 5, 10}. Data is N2 = 500 gaussian random data with $\mu = [0, 1.5]$ and $\Sigma = [0.9, 0.4 ; 0.4 , 0.9]$.

# Figure 1-4: Part 2

## Figure 1

**Figure 2**

**Figure 3**

**Figure 4**

## Figure 5-8: Part 3

### Figure 5

**Figure 6**

**Figure 7**

**Figure 8**

**Figure 9-12: Part 4 - $N_1$**

**Figure 10**

**Figure 11**

**Figure 12**

**Figure 13 – 16: Part 4 – N$_2$**

**Figure 13**

**Figure 14**

**Figure 15**

**Figure 16**