# Assignment 2, COMP4702

Roy Portas, 43560846

May 2, 2017

# Question 4.2

```matlab
1  % Q2
2
3  a = randn(200, 2);
4  b = a + 4;
5  c = a;
6
7  c(:, 1) = 3 * c(:, 1);
8  c = c - 4;
9
10 d = [a; b];
11 e = [a; b; c];
12
13 % hold on;
14 % plot(a(:, 1), a(:, 2), '+');
15 % plot(b(:, 1), b(:, 2), 'o');
16 % plot(c(:, 1), c(:, 2), '*');
17
18 % Use the first dataset
19 data = e;
20
21
22 figure;
23 subplot(3, 2, 1);
24 hold on;
25 % ksdensity(data, 'PlotFcn', 'contour');
26 plot(a(:, 1), a(:, 2), '+');
27 plot(b(:, 1), b(:, 2), 'o');
28 plot(c(:, 1), c(:, 2), '*');
29
30 title('Contours Overlay');
31
32 % Calculate the grid
33 [f, xi] = ksdensity(data); x = linspace(min(xi(:, 1)),max(xi(:, 1)
     ));
34 y = linspace(min(xi(:, 2)),max(xi(:, 2)));
35 [xq,yq] = meshgrid(x,y);
36 z = griddata(xi(:, 1),xi(:, 2),f,xq,yq);
37
38 % We now have x, y, z that can be used to get the gradient at any
     point
39
40 contour(x, y, z);
```

```matlab
41  xlim([-10, 10]);
42  ylim([-10, 10]);
43  hold off;
44
45  copy = data;
46  new_points = copy;
47
48  for i = 1:5
49      new_points = step(new_points, x, y, z);
50      subplot(3, 2, i + 1);
51      hold on;
52      % ksdensity(data, 'PlotFcn', 'contour');
53      scatter(new_points(:, 1), new_points(:, 2));
54      title(sprintf('Step %d', i));
55      xlim([-10, 10]);
56      ylim([-10, 10]);
57      hold off;
58
59  end
60
61
62  function dist = euclid_distance(point1, point2)
63      dist = sqrt((point1(1) - point2(1))^2 + (point1(2) - point2(2)
            )^2);
64  end
65
66  function new_points = step(points, x, y, z)
67      % Calibration factor (lambda)
68      max_distance = 1.8;
69
70      new_points = zeros(length(points), 2);
71      for i = 1:length(points)
72          point = points(i, :);
73
74          numerator = 0;
75          denominator = 0;
76
77          for j = 1:length(points)
78              other_point = points(j, :);
79              if euclid_distance(point, other_point) < max_distance
80
81                  % Calculate part of sum
82                  distance = euclid_distance(point, other_point);
83                  numerator = numerator + (distance * other_point);
84                  denominator = denominator + distance;
```

```
85              end
86          end
87
88          mx = numerator / denominator;
89          new_points(i, :) = mx(1, :);
90      end
91  end
```
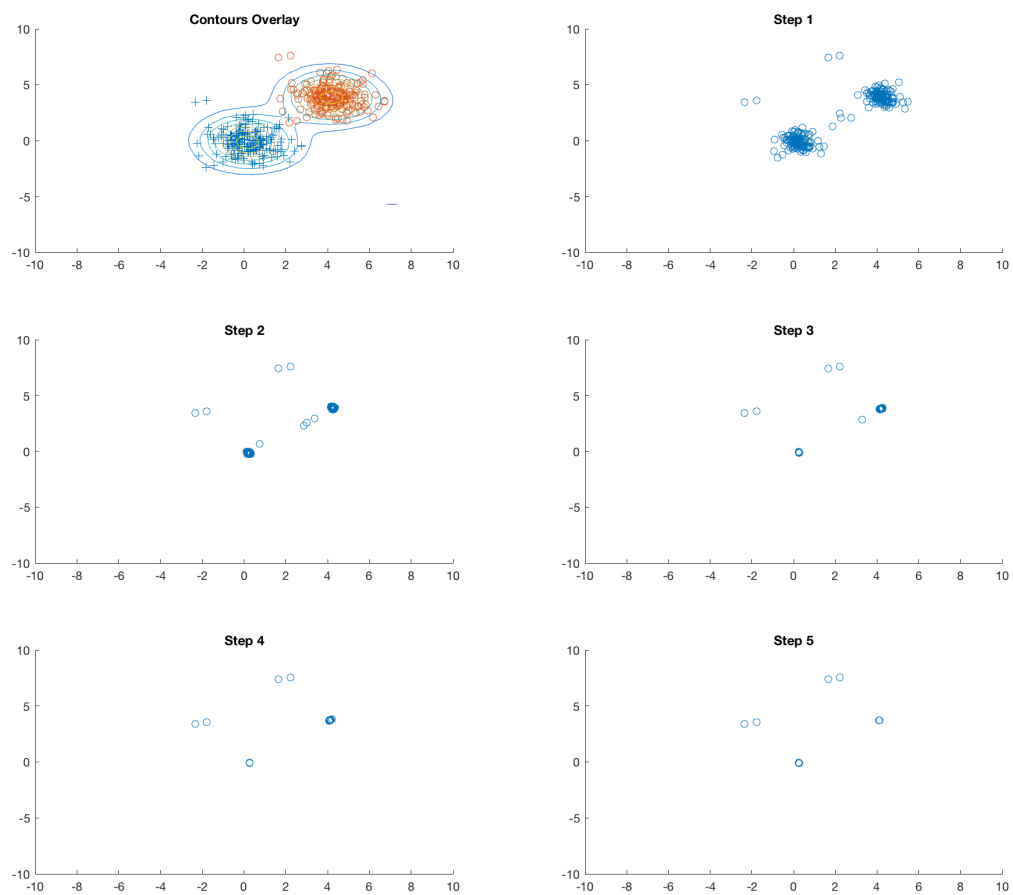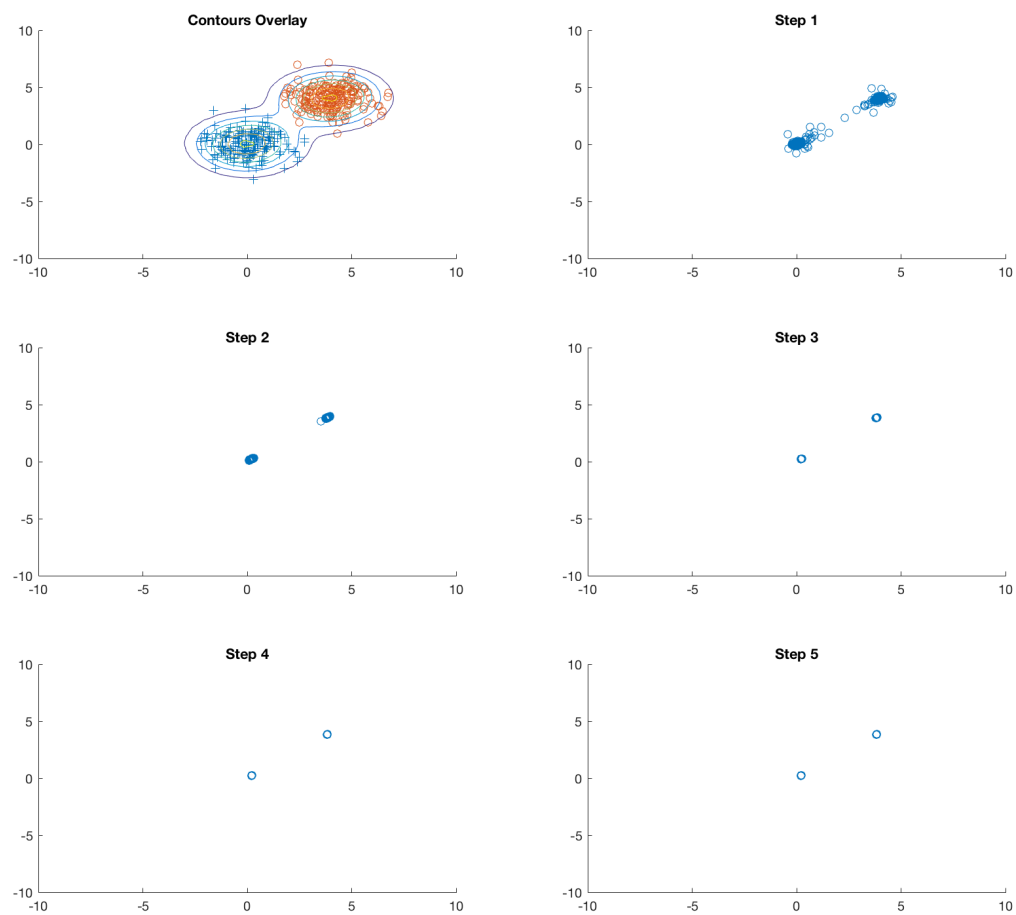
# Question 4.3



Figure 1: 2 Classes, Lambda = 1.8
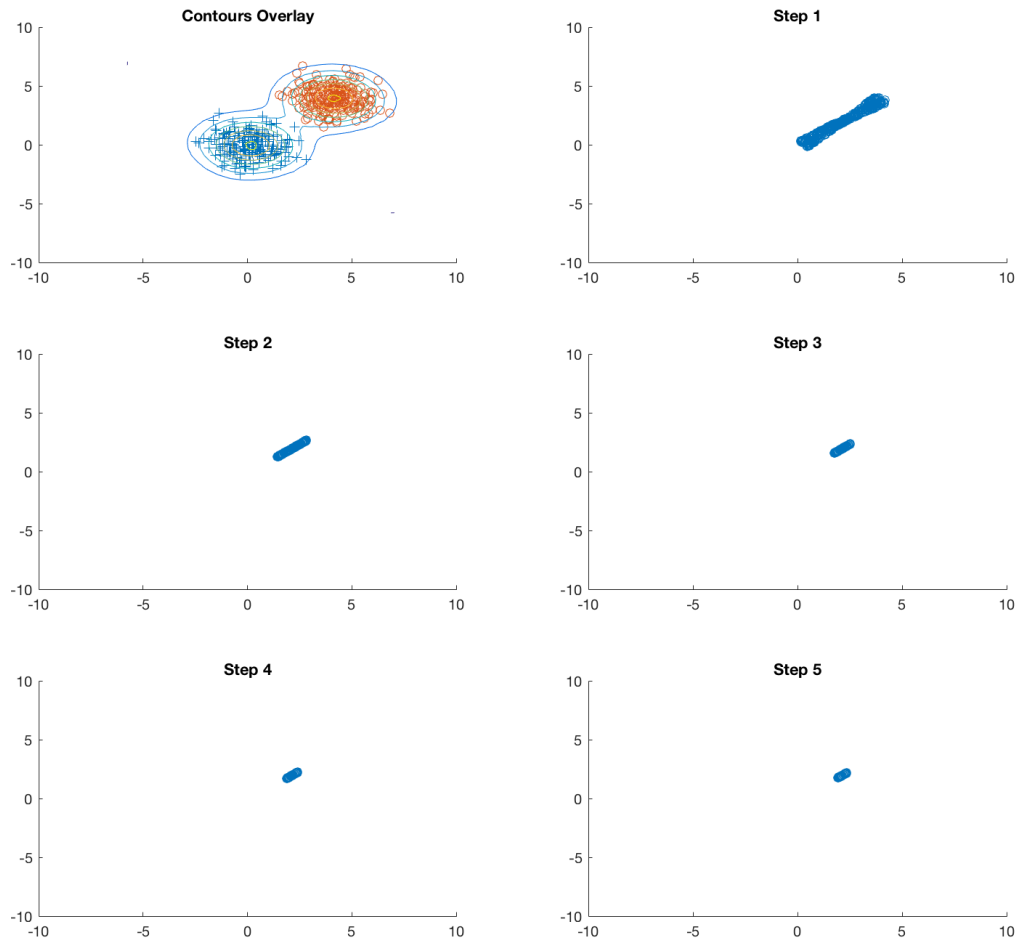
Figure 2: 2 Classes, Lambda = 3

Figure 3: 2 Classes, Lambda = 5

For the 2 class problem, a lambda value set to 3 provided the best result. A lambda of 1.8 cause a few outliers not to shift towards the mean value. Whereas a lambda value of 5 caused all the points to shift towards one of the means, which is not desired. A lambda value of 3 shifted all the points to the two mean values without leaving outliers, thus is the best lambda value out of the three.
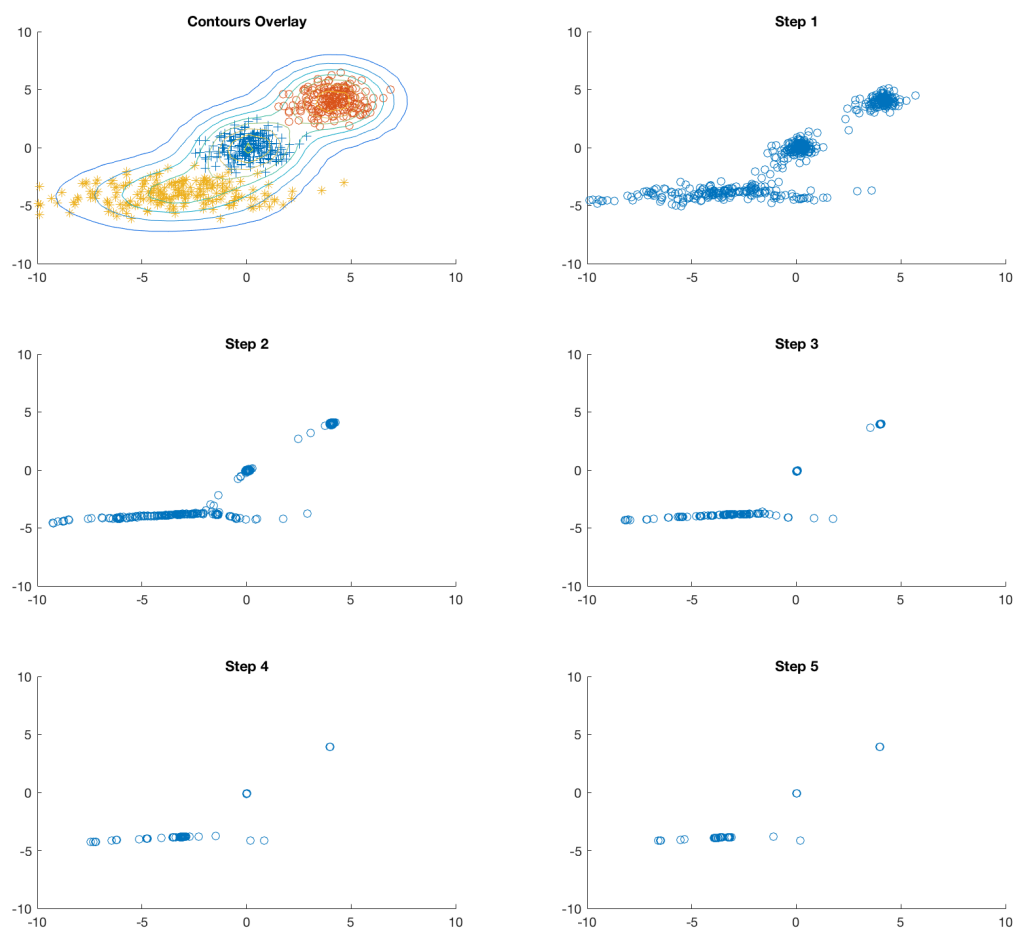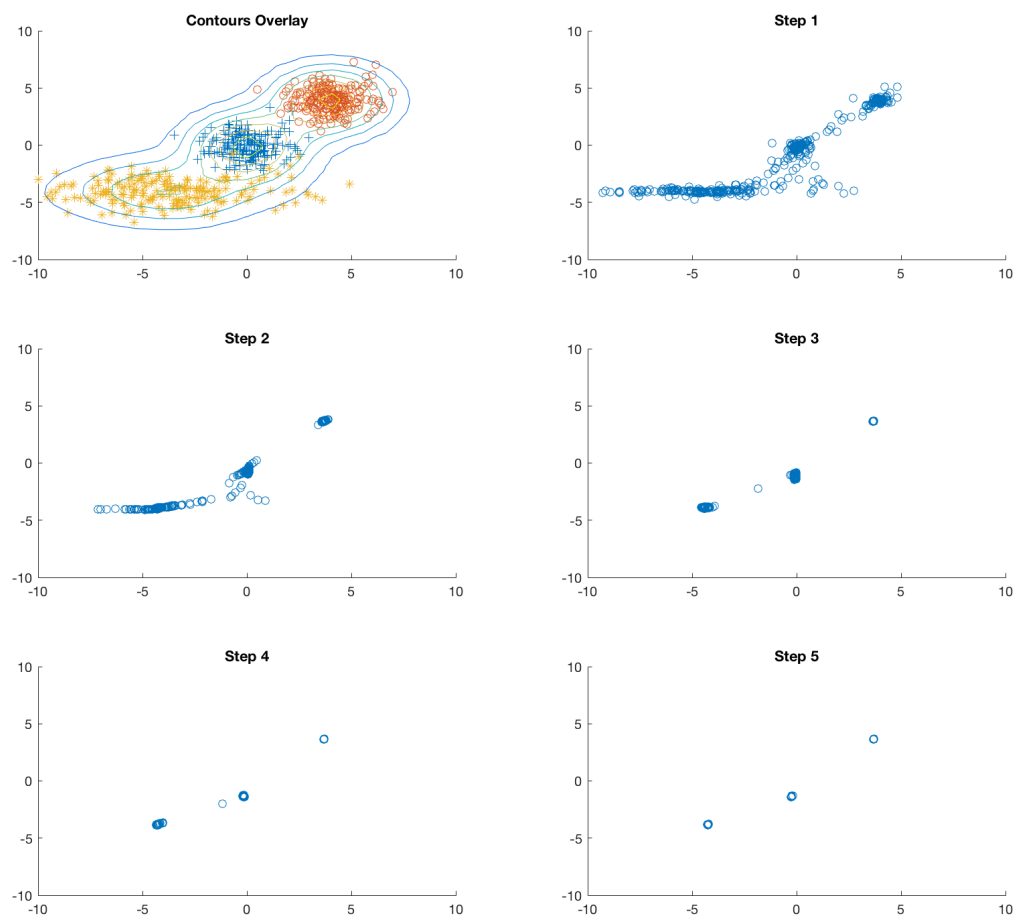
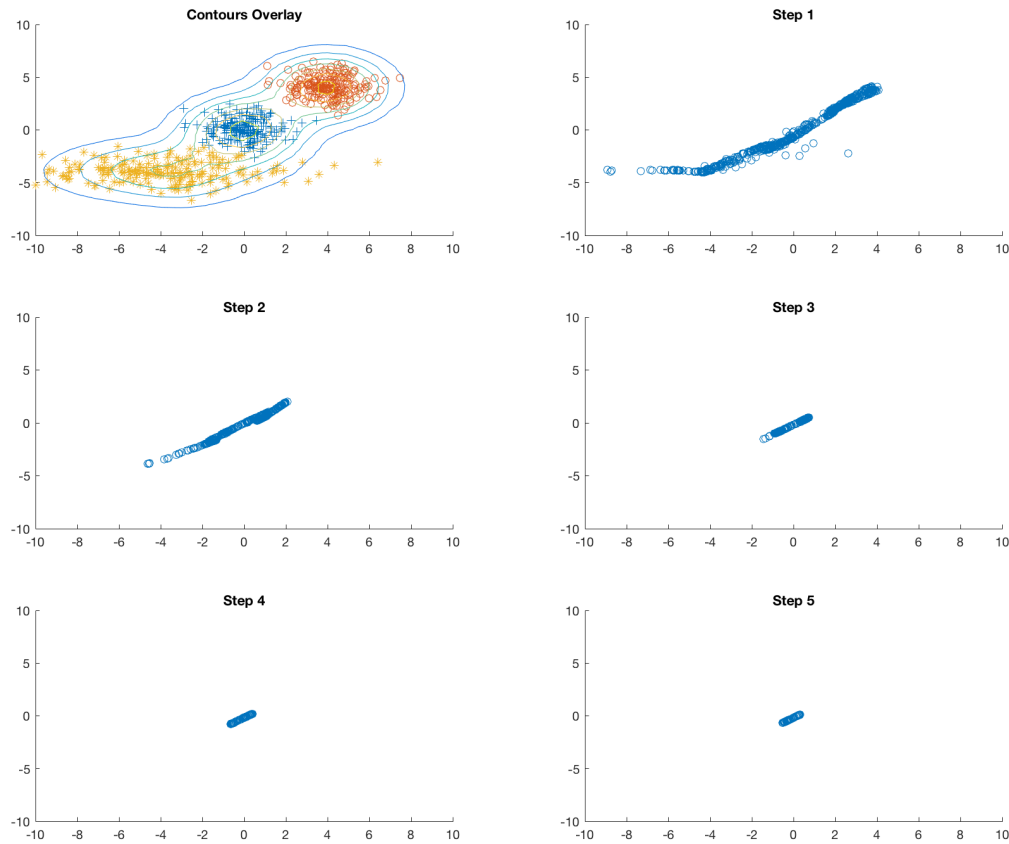Figure 4: 3 Classes, Lambda = 1.8

Figure 5: 3 Classes, Lambda = 3

Figure 6: 3 Classes, Lambda = 5

The 3 class problem showed similar results with the same lambda values 1.8, 3 and 5. Again the 1.8 lambda value failed to shift some outliers towards the mean and the 5 lambda value shifted all of them towards a single point. The lambda value 3 correctly shifted all the points to their respective means, thus the lambda value of 3 was the best choice out of the three numbers.

# Question 5.1

```matlab
%
% Principle Component Analysis
%

function result = pca(data)
    m = mean(data);
```

```matlab
7      S = cov(data − m);
8       [evec, eval] = eigs(S);

10      % Sort the eigenvalues
11       [y, i] = sort(diag(eval), 'descend');
12      % Sort the eigenvectors columns by the eigenvalue indexes
13       evec = evec(:, i);

15      % PCA only works if you subtract the mean
16       result = evec' * (data − m)';
17   end
```
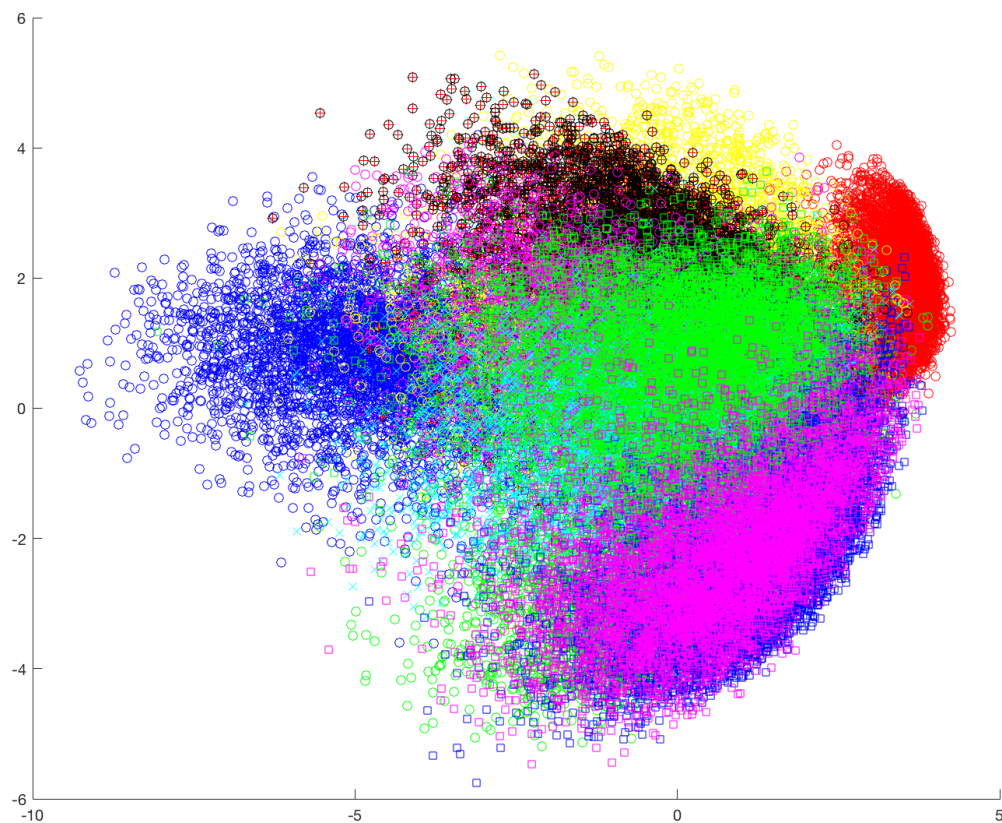
# Question 5.2

**a**



Figure 7: PCA on MNIST dataset

# b

The first principle component accounts for 5.116% of the data, whereas the second principle component accounts for 3.7414% of the data.
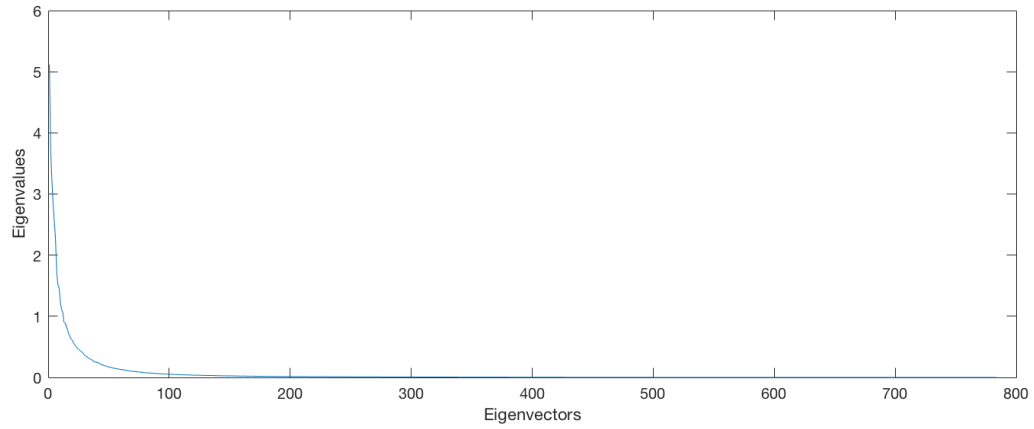
# c



Figure 8: Scree Graph

# Question 5.6


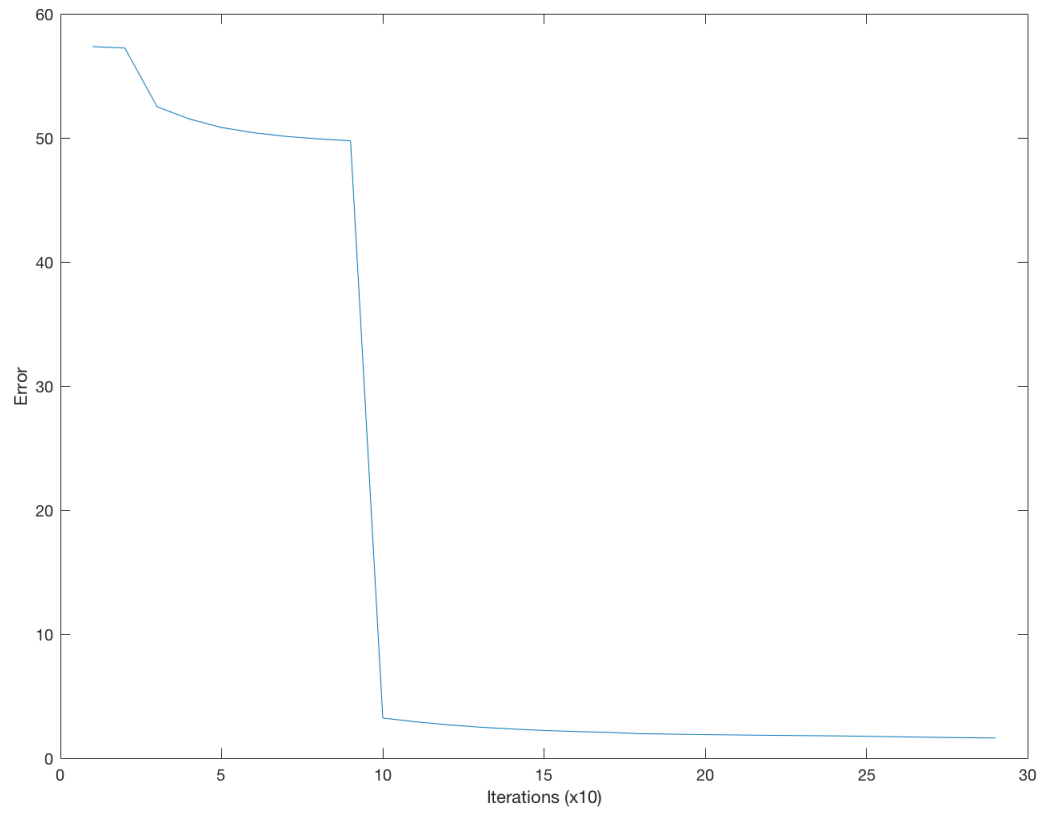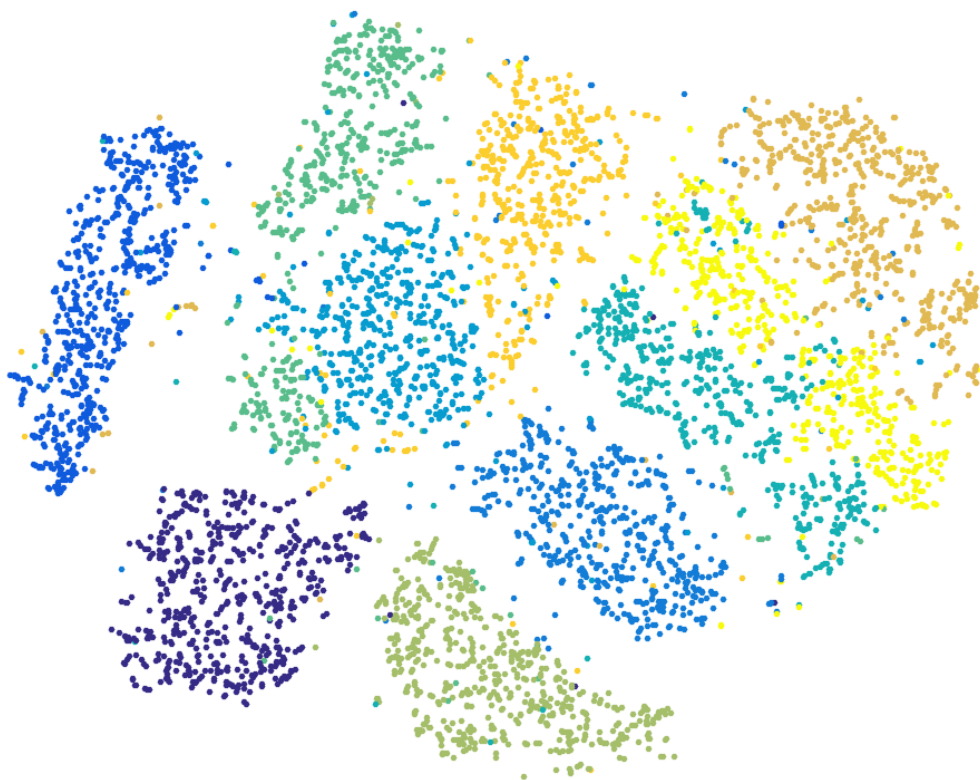
Figure 9: Error vs Iteration

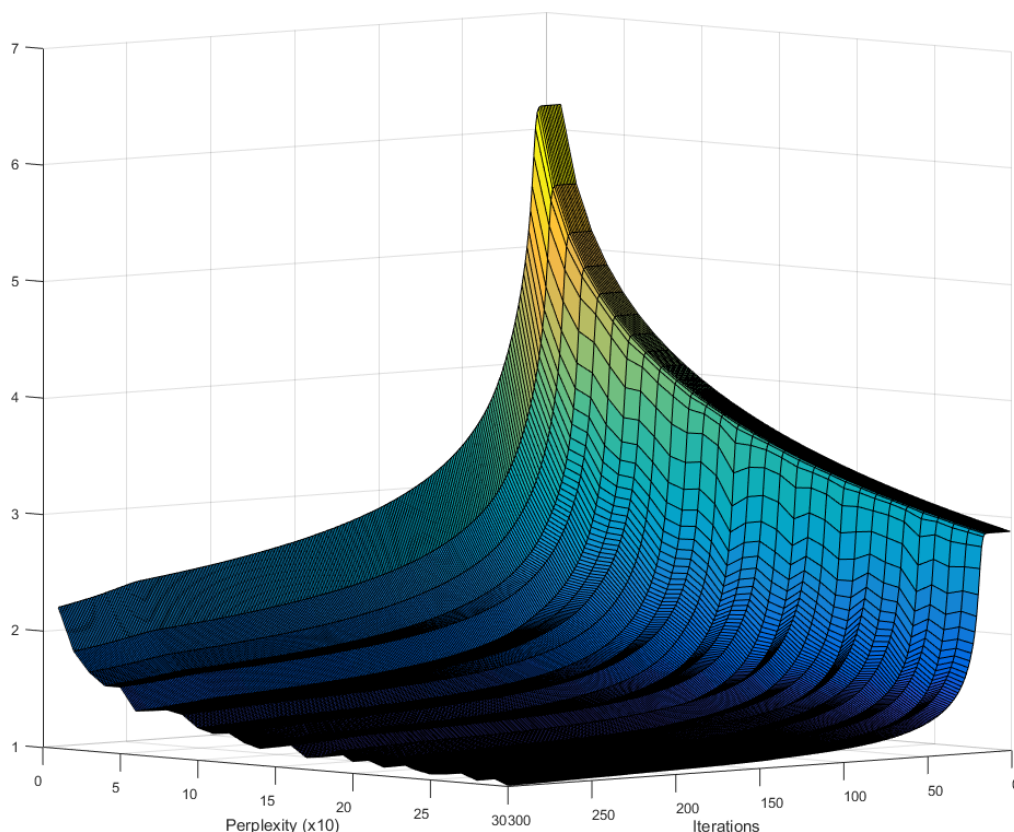Figure 10: Iteration 300

# Question 5.8



Figure 11: Perplexity and Iterations on Cost

The perplexity determines how to balance the attention between local and global aspects of the dataset. The graph shows the error change on iterations, thus we should choose a perplexity the reduces the error as quickly as possible, such that the algorithm requires less iterations to find a good solution. A simple heuristic to choose a good perplexity would be to look for the greatest rate of change in the error, by inspecting the chart we can see that the lower perplexities have far greater rate of change. Further inspection showed that a perplexity value of 2 provided the best rate of change.

This chart differs from the one in figure 9, as the other chart as a large dip in the error at the 100 iteration mark. This is because of the lines that were commented out, as those lines introduce a larger momentum that decays over each iteration. Line 87-89 of `tsne_p.m` tells t-SNE to stop lying about the p-values when it reaches the 100th iteration, which explains the massive drop in error.
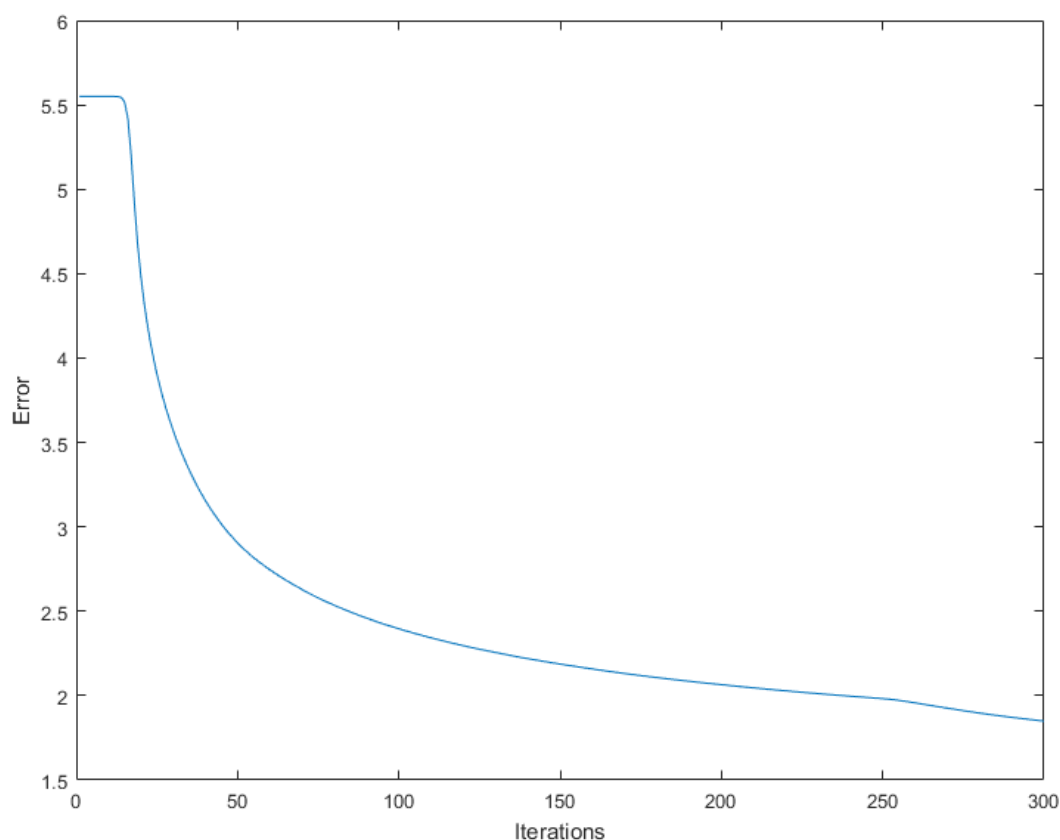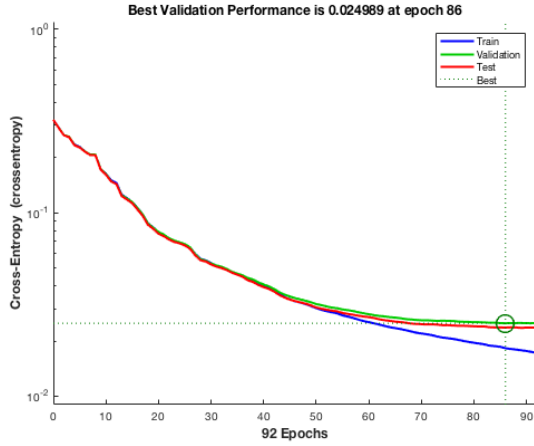
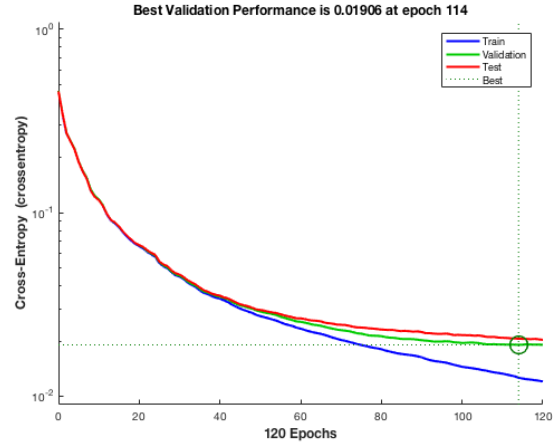Figure 12: Error on Perplexity

# Question 6.3

Using the code generated from the `nprtool`, the MNIST dataset can be easily loaded however there are various design choices that should be addressed.

The first design choice is the network training function. I chose the `trainscg` function as it does not require a line search on each iteration, which is described as a computationally intensive operation in the Matlab documentation. The Matlab describes the downside of adopting this function is it takes more iterations to converge than the other algorithms, but the speed increase is worth it.
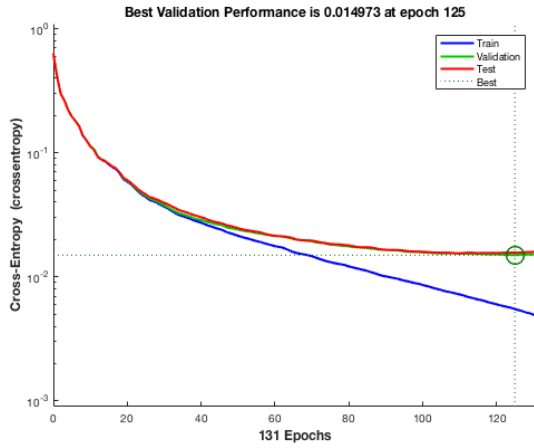
Another design parameter is the number of neurons in the hidden layers. The documentation for Scikit learn states the general rule of thumb for choosing values for the hidden layers is a value between the number of neurons in the input and output layers. As we are dealing with 784 input layers and 10 output layers, any value between those would work. To choose a suitable value, the network was trained against 15, 30, 45 and 60 hidden neurons,
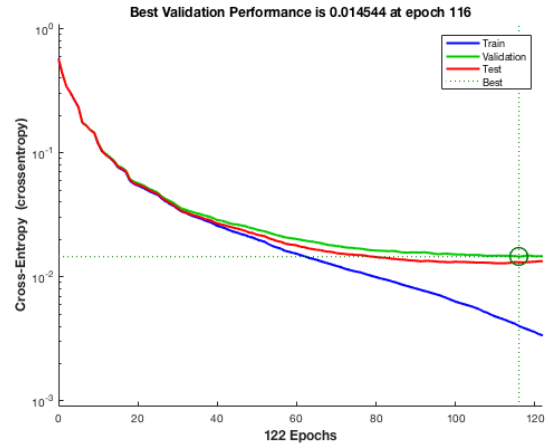
(a) 15 hidden neurons training



(b) 30 hidden neurons training



(c) 45 hidden neurons training



(d) 60 hidden neurons training

below is a chart showing the results.

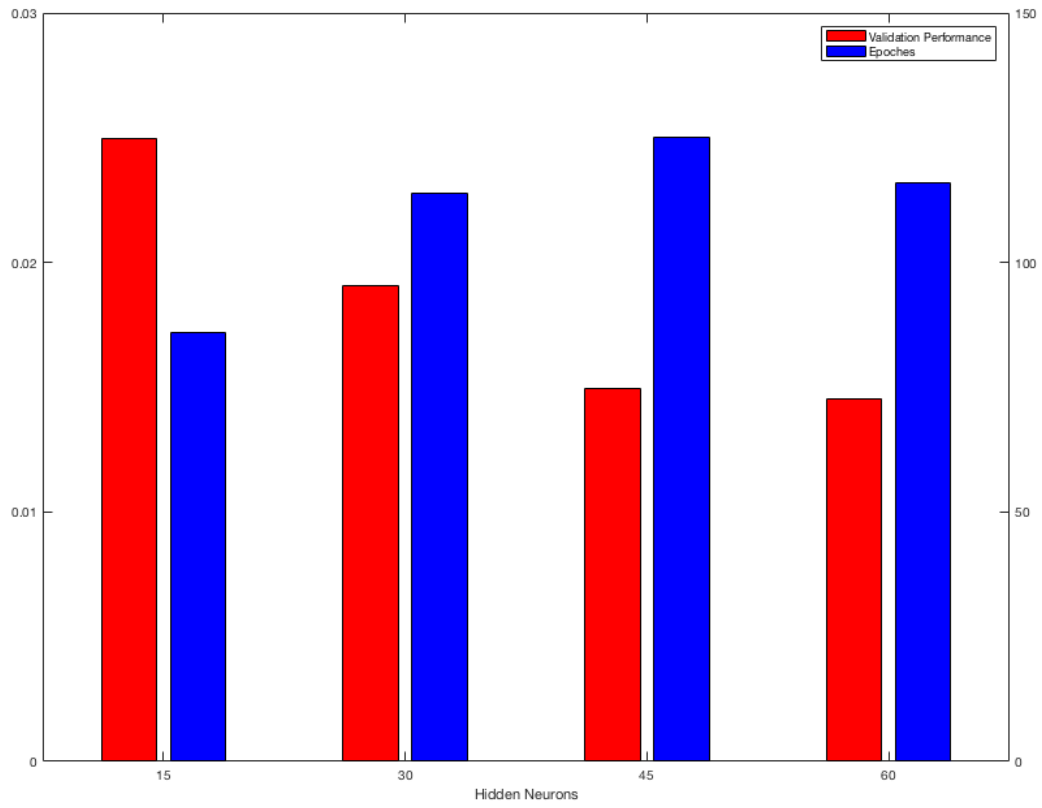The results of the training was summarised into the following graph.

Figure 14: Number of Neurons vs Performance and Epochs

It can be seen that after 45 neurons, there are only marginal increases in performance, but the number of epochs until converging increases, thus increasing the running time of the program. Thus 60 neurons was selected.
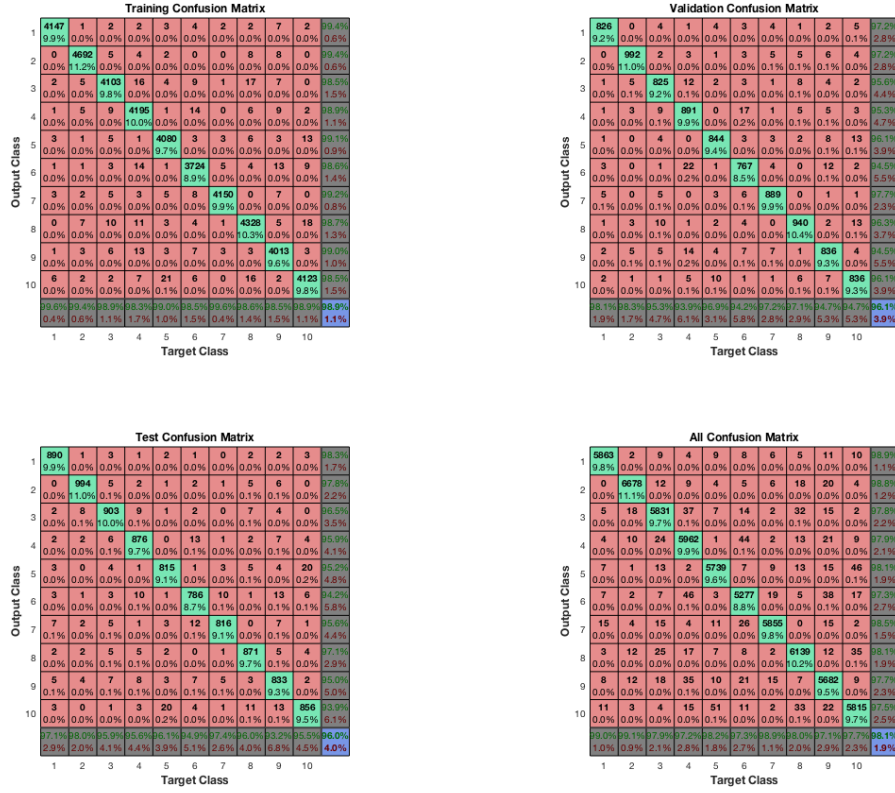
Figure 15: 60 hidden neurons confusion matrix

The above confusion matrix chart shows the error in the algorithm. The validation confusion matrix shows that the overall error is 3.9%. This can be compared to Figure 9 in the paper "Gradient-Based Learning Applied to Document Recognition", which found the error rate to be 4.7%.