

## Prac 4

## Q2

**Process:** blurring**Code Implementation:** Function**Inputs:** x: data. lambda:  $\lambda$  (radius)

```

1  function prac4_test(x,lambda)
2  %-----Plot X START-----
3  a = randn(200,2);
4  b=a+4;
5  c=a;
6  c(:,1) = 3*c(:,1);
7  c=c-4;
8  d=[a;b];
9  e=[a;b;c];
10 plot(a(:,1),a(:,2),'+');
11 hold on
12 plot(b(:,1),b(:,2),'o');
13 plot(c(:,1),c(:,2),'*');
14 %-----Plot X END-----
15 times = 0;
16 ct = 0;
17 old_Center = 0;
18 this_Center = 0;
19 meanlist = [];
20
21 for i = 1:size(x,1)
22     this_Center = x(i,:);
23     times = times + 1
24
25     while(this_Center ~= old_Center)
26         index = 0;
27         s = 0;
28         for j = 1:size(x,1)
29             o = x(j,:);
30             k = norm(o-this_Center);
31             if k <= lambda
32                 index = index+1;
33                 s = s + o;
34             end
35             s = s + o;
36         end
37         old_Center = this_Center;
38         this_Center = s ./ index;
39         if isnan(this_Center)
40             break;
41         end
42     end
43
44     if (~ismember(this_Center,meanlist))
45         ct = ct+1;
46         meanlist(ct,:) = this_Center;
47         plot(this_Center(:,1),this_Center(:,2),'kd','MarkerSize',8,'MarkerFaceColor',[.49 1 .63])
48         hold on;
49     end
50 end
51 hold off;
52
53 end

```

### Q3

**Cluster Centres:** green diamond mark

**Red Marks ('o'):** b

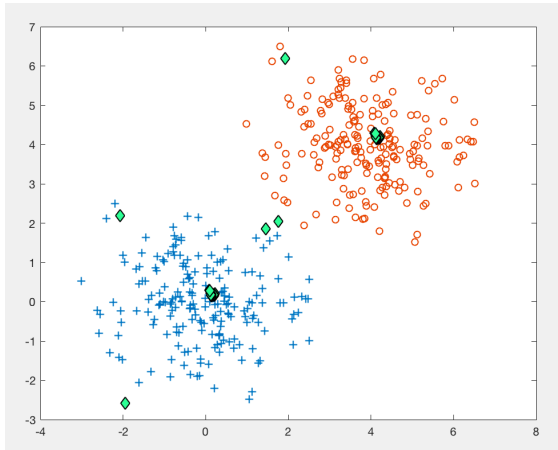
**Blue Marks ('+'):** a

**Yellow Marks ('\*'):** c

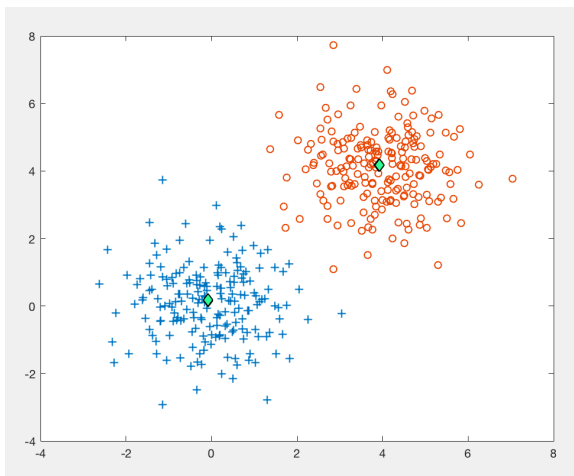
**Test Data with 3 different lambda**(*Blue region is not included in this test*):

$d=[a;b];$

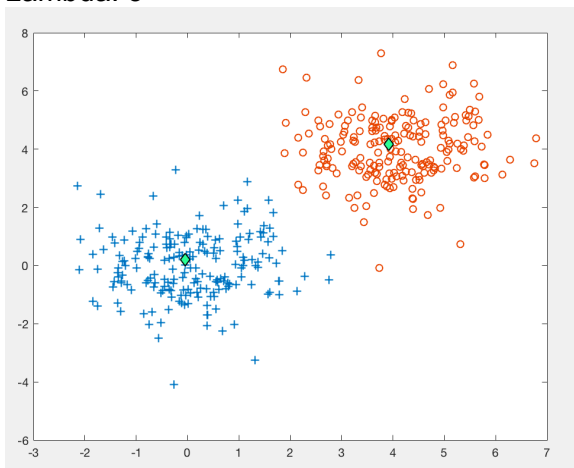
Lambda: 1



Lambda: 2



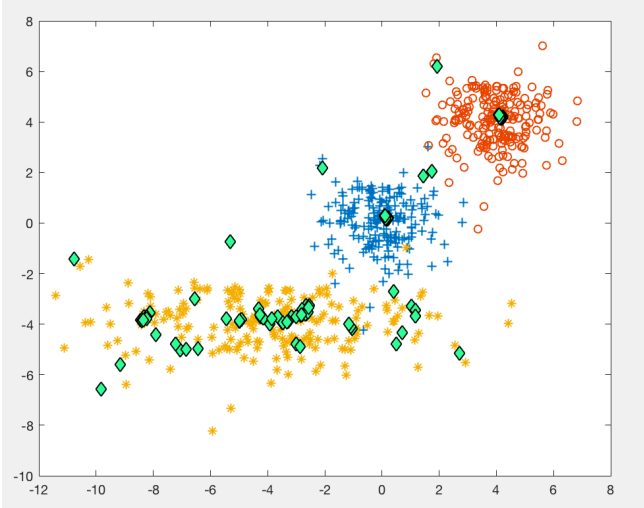
Lambda: 3



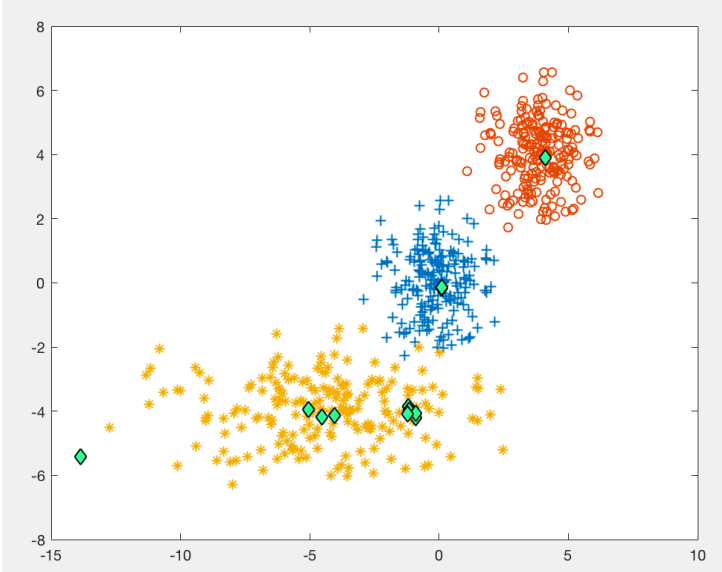
Test Data with 3 different lambda:

e=[a;b;c];

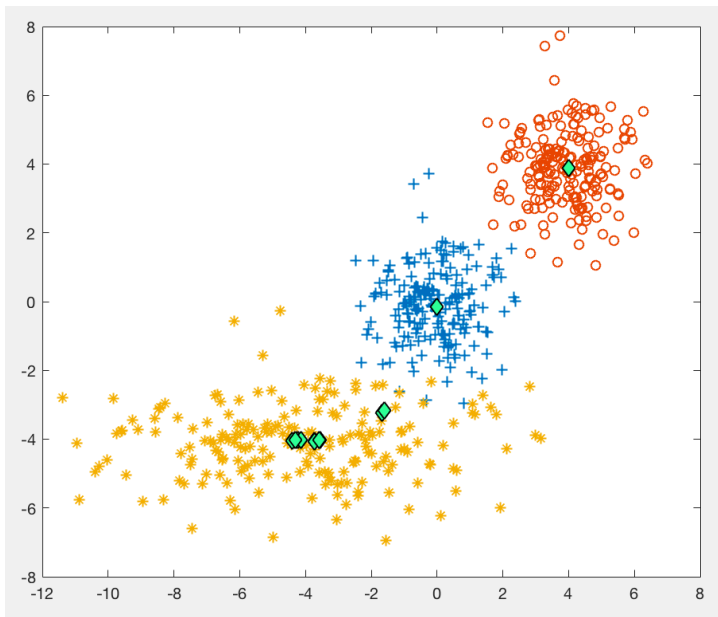
Lambda: 1



Lambda: 2



Lambda: 3



Larger lambda can get much more concentrated and less cluster centres. This effect is apparent on yellow regions.

Because larger lambda can make bigger move step to next mean centre and less iteration, thus larger lambda also affect the speed of running the mean shift algorithm.

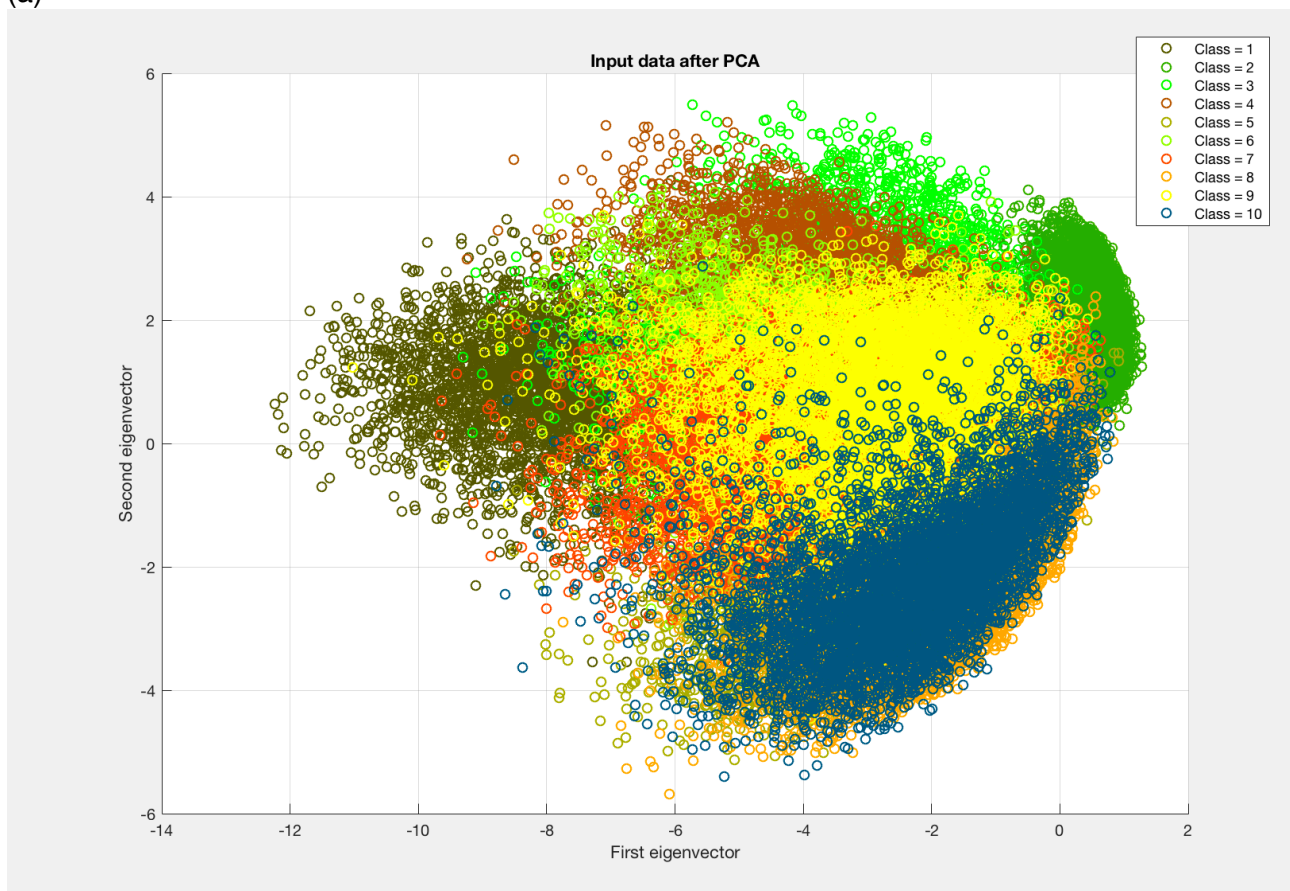
# Prac 5

## Q1

```
1 function [povResult,reducedData] = PCA(train_X,train_labels,n)
2 %prac5_1 Principal Component Analysis
3 %
4 % [P, beta] = PCA(dataset1, dataset2, number)
5 % Inputs:
6 % dataset1: Raw data (For MNIST: train_X)
7 % dataset2: Data labels (For MNIST: train_labels)
8 % number: Expected classes (But only Max 64 can be displayed on
9 %         Figure 1)
10 % Outputs:
11 % ReducedData: Reduced to 2D, included First & Second principal component
12 % povResult: Proportion of variance for Q2(b)
13 %% Q1
14 x = train_X; %n = 10;
15 data = cov(x);
16 [vector,value] = eig(data);
17 %First principal component: 1st largest -> Dimention Reduction
18 X1 = x*vector(:,end);
19 %Second principal component: 2nd largest -> Dimention Reduction
20 X2 = x*vector(:,end-1);
21 %Reduced Data
22 reducedData = [X1,X2];
23 size(reducedData)
24 %Added lables onto reduced data
25 X1 = horzcat(X1,train_labels);
26 X2 = horzcat(X2,train_labels);
27 %% Q2
28 c=colormap;
29 eval0 = eig(data); % inputed_raw eigenvalues
30 labels = [];
31 for ii = 1:n
32     a = X1(X1(:,2)==ii);%x_axis
33     b = X2(X2(:,2)==ii);%y_axis
34     scatter(a,b,[],c(ii,:),:));
35     legendInfo{ii} = ['Class = ' num2str(ii)];
36     hold on;
37 end
38
39 grid on;
40 legend(legendInfo);
41 xlabel('First eigenvector');
42 ylabel('Second eigenvector');
43 title('Input data after PCA');
44 %% Q3 & Q4
45 sort_value = sort(eval0,'descend');
46 figure,subplot(2,1,1);
47 plot(sort_value,'+-');
48 grid on;
49 xlabel('Eigenvectors')
50 ylabel('Eigenvalues')
51 title('(a) Scree graph for Input data');
52
53 hold on;
54 %Proportion of variance: used feature(1->k)/inputed_raw features(1->d)
55 povResult = (eval0(end)+ eval0(end-1))/sum(eval0);
56 subplot(2,1,2);
57 cumsum_sort_value = cumsum(sort_value);
58 plot(cumsum_sort_value/sum(eval0),'+-');
59 hold on;
60 txt = ['Q2\b = ',num2str(povResult)];
61 plot(povResult,'kd','MarkerSize',8,'MarkerFaceColor',[.49 1 .63])
62 hoz_line = reffline([0 povResult]);
63 hoz_line.Color = 'r';
64 text(0,povResult,txt,'HorizontalAligment','left')
65 grid on;
66 xlabel('Eigenvectors')
67 ylabel('Prop. of var.')
68 title('(b) Proportion of variance explained');
69 hold off;
70 end
```

Q2

(a)



(b)

Percentage of the data variance is accounted for by the first two principal components:

0.16801 (Green diamond mark with red line on below graph)

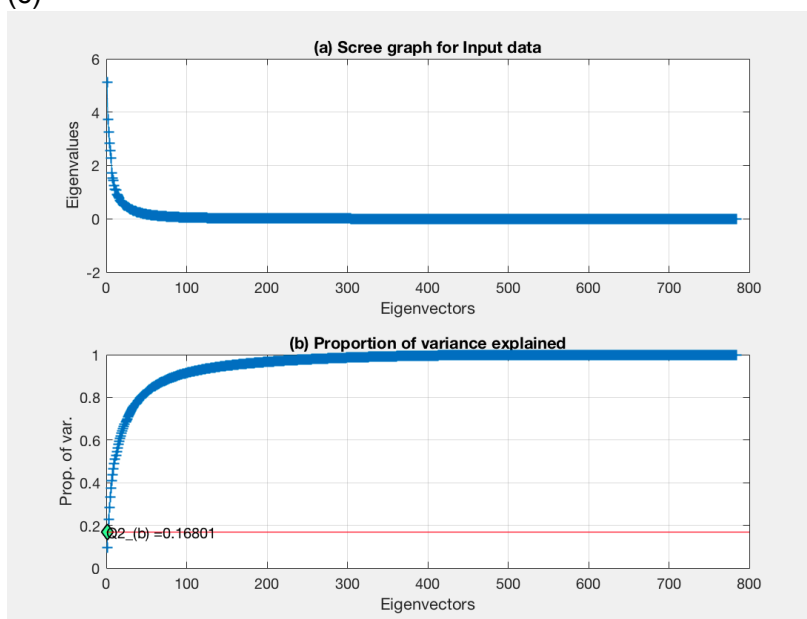
Code: `povResult = (eval0(end)+ eval0(end-1))/sum(eval0);`

The **eval0(end)** means the first principal component.

The **eval0(end-1)** means the second principal component.

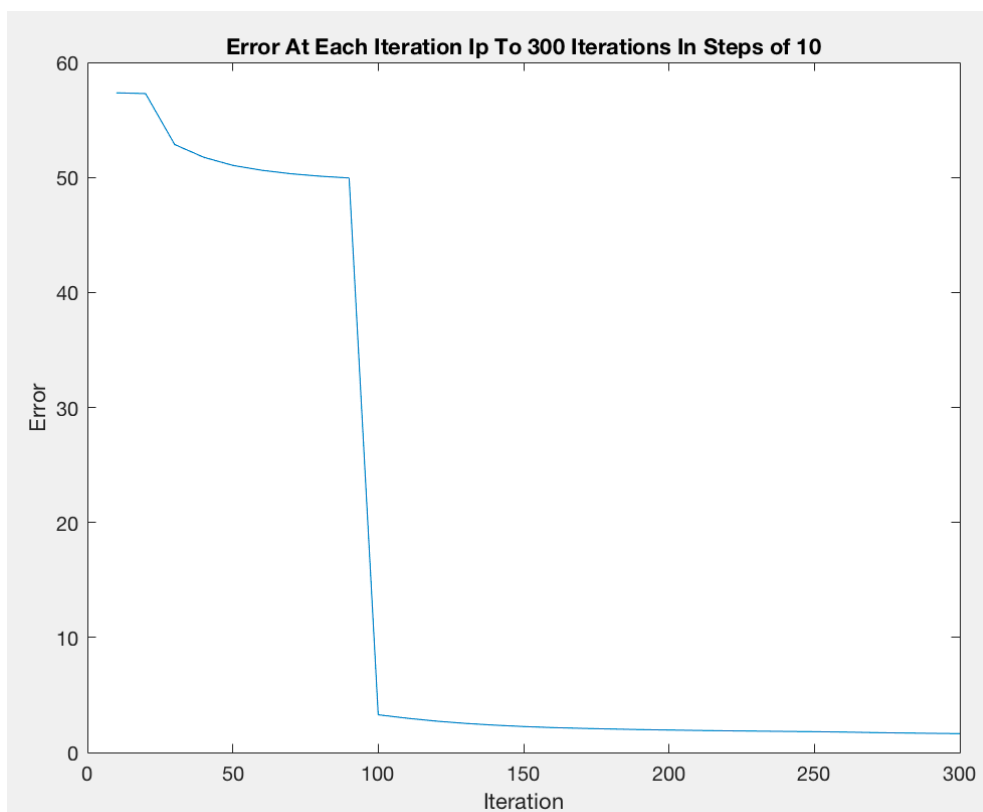
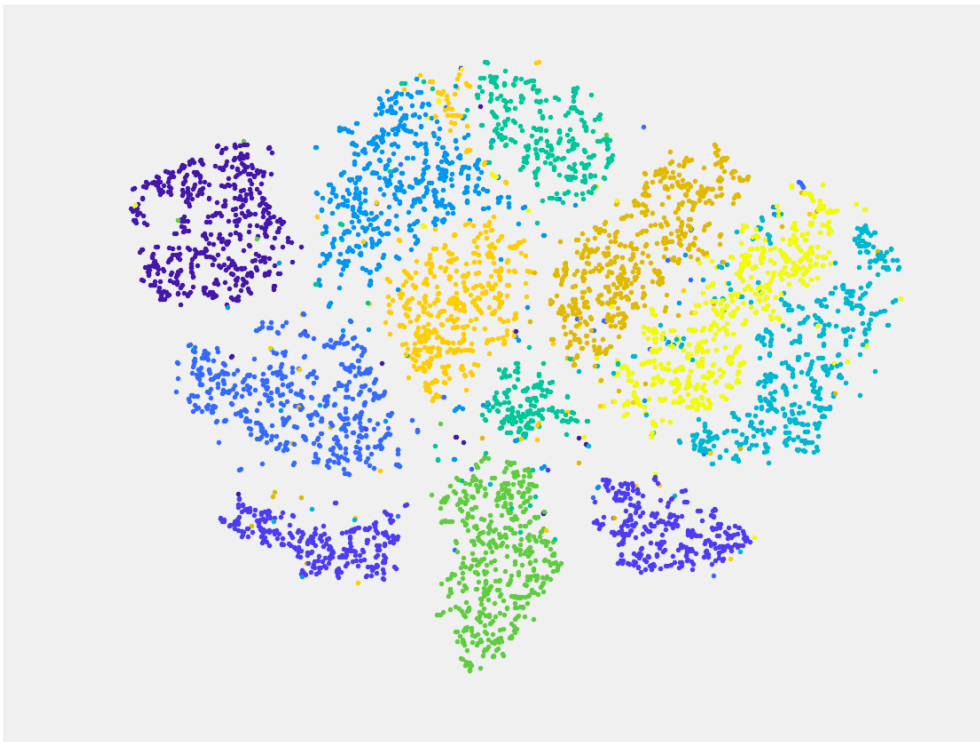
The **sum(eval0)** means the total Eigenvalues got from the data.

(c)



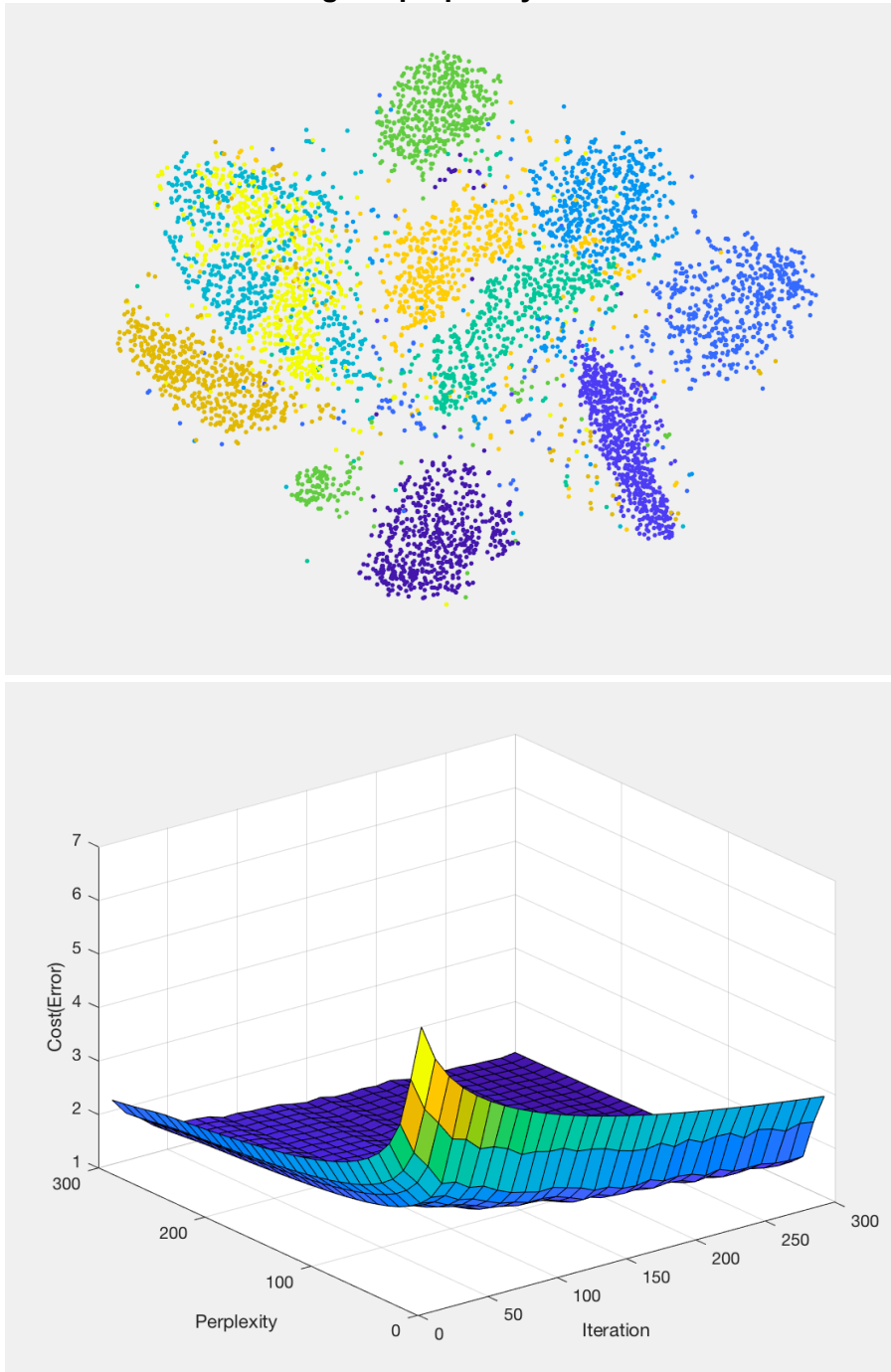
Q6

Screenshot after 300 iterations:



Q8

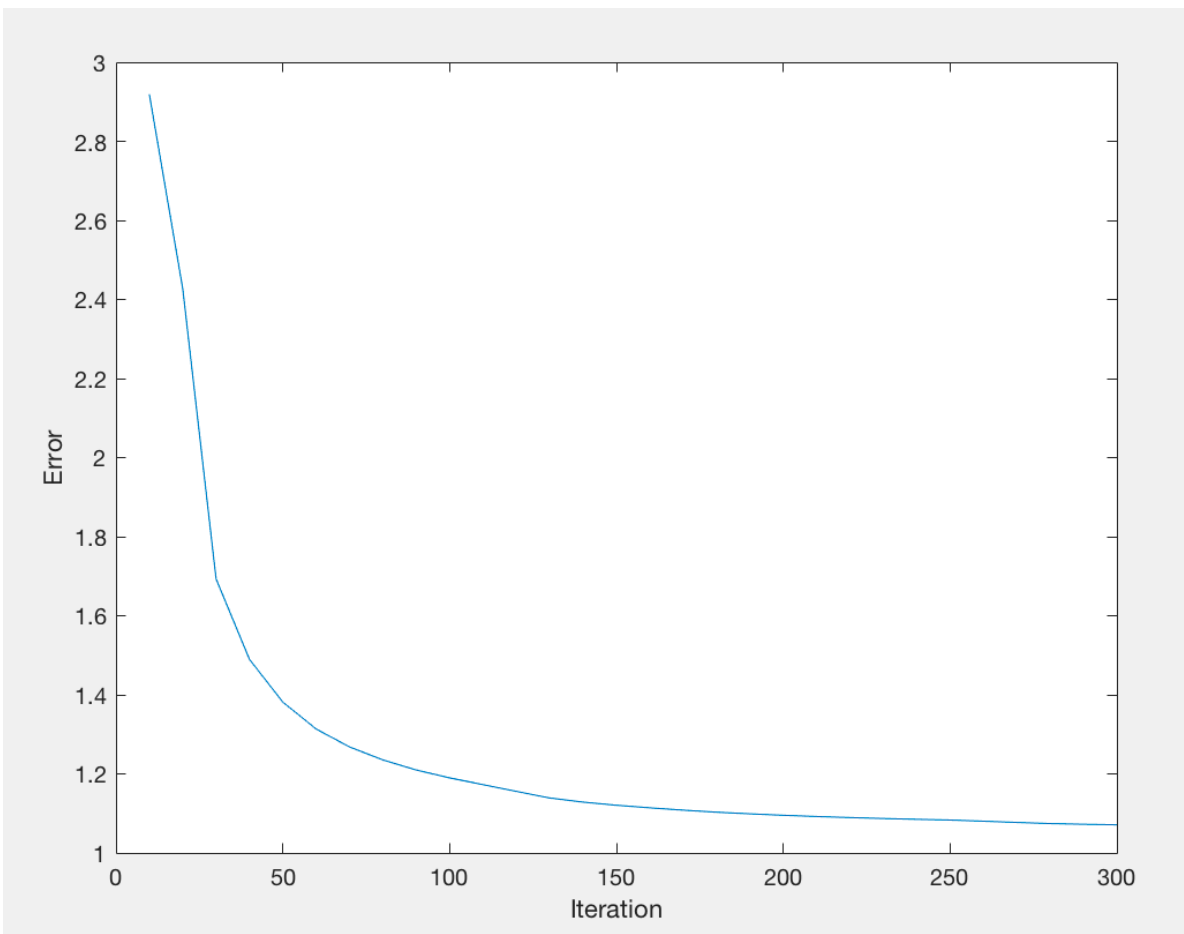
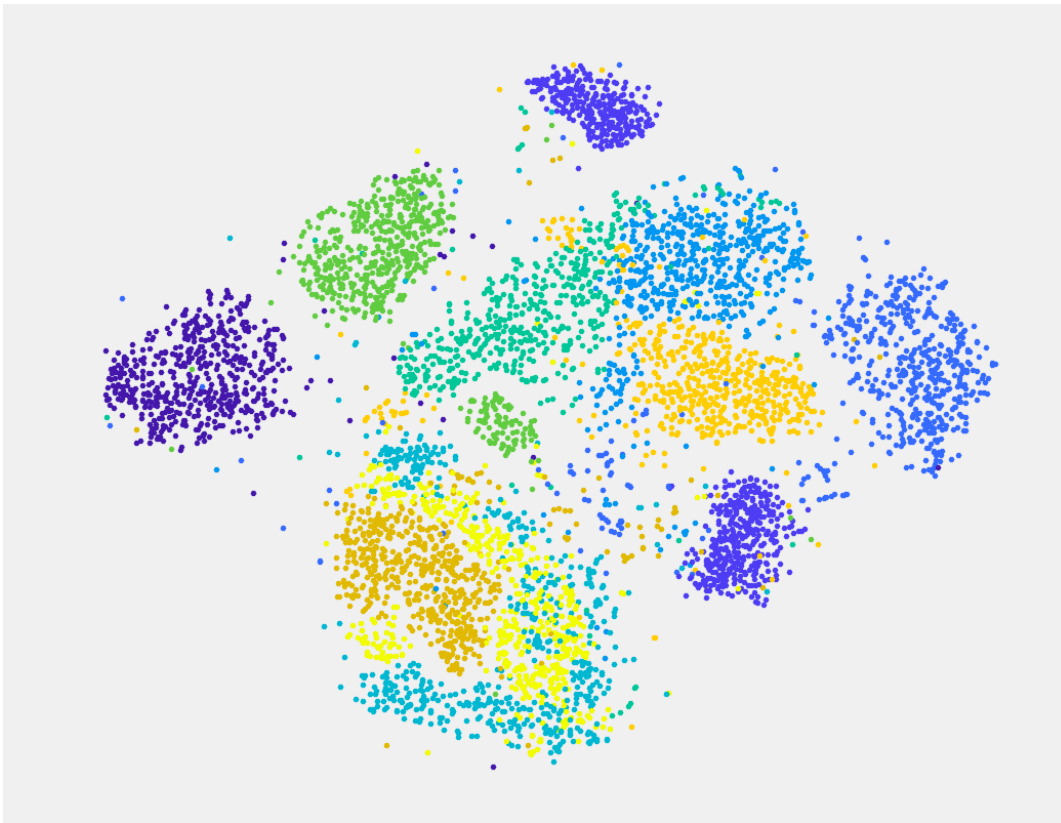
Screenshot after setting the perplexity



Heuristic for choosing the perplexity:

When watch at the 3D graph, when Cost becomes lower, the perplexity becomes larger, thus, the best choice of perplexity by now will be 290, and it is smaller than 300, and it uses less time to go through all possible perplexity.





The 2D visualisation is not good as Q6, Q6 got clearer classification. Perplexity is also bigger than Q6, thus perplexity really affect the classification results.

Because there are some codes deleted, search area becomes smaller, and the errors are also much smaller than Q6.

## Prac 6

### Q3

*gd: Gradient descent backpropagation*

*scg: Scaled conjugate gradient backpropagation*

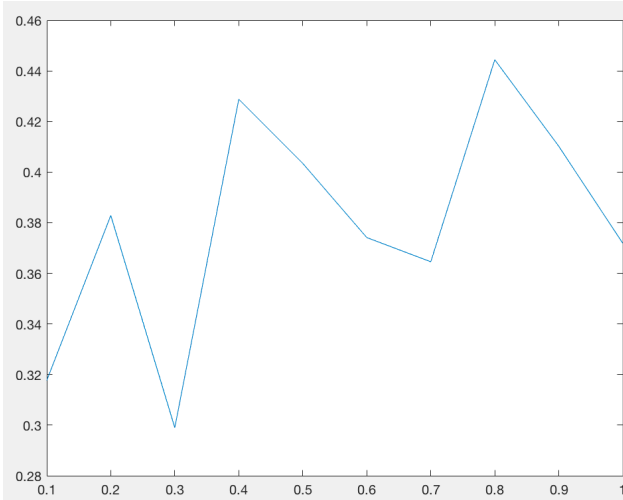
Tested learning rate, hidden layer size, epochs independently.

**Learning Rate Range(x\_axis): 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1**

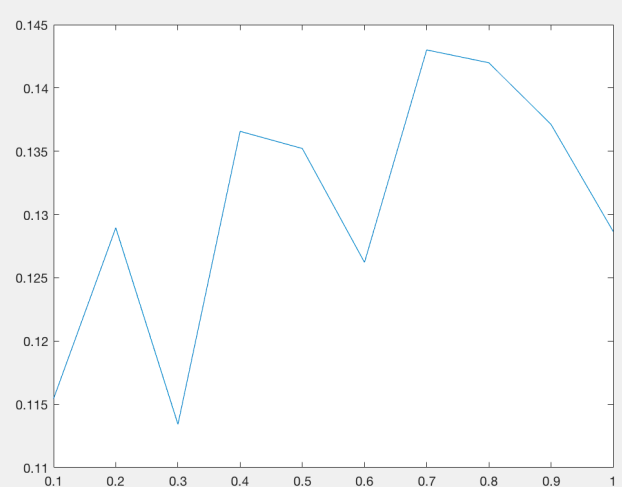
*<Learning rate can only be positive, and the maximum is 1, default is 0.01>*

**Training Function: gd**

Left Graph: y\_axis: percentage Error

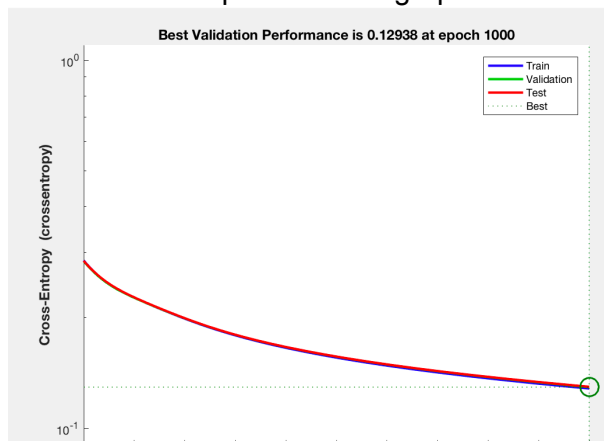


Right Graph: y\_axis: performance



After plotting both results of percentage errors and performance values for each learning rate. The best learning rate in my test samples is 0.3, because both graph and training functions indicate the local minimum point at 0.3. The graphs also oscillate apparently at each 0.2 difference.

gd performance during running each learning rate shows: there is no overfitting or under fitting. And most of the performance graphs are similar to below graph.

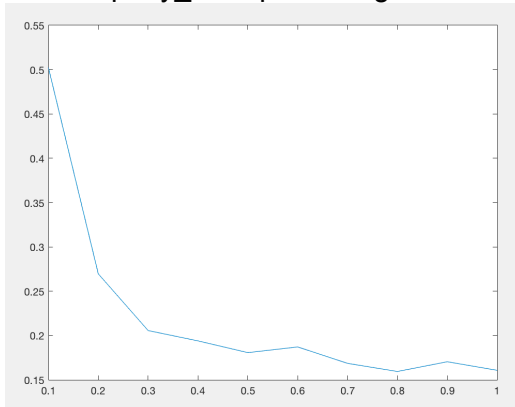


*<Because only gradient descent algorithm has learning rate, thus I only tested the gd.>*

## Hidden Layer Size Range(x\_axis): 10 20 30 40 50 60 70 80 90 100

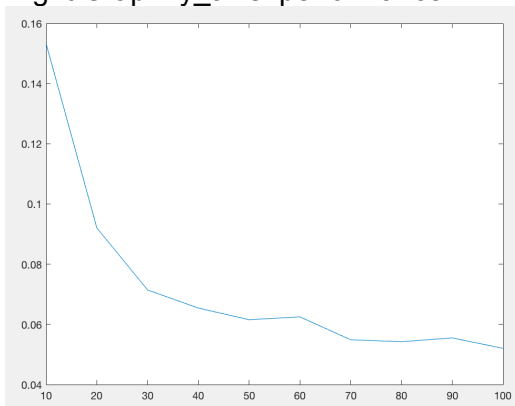
Training Function: gd

Left Graph: y\_axis: percentage Error



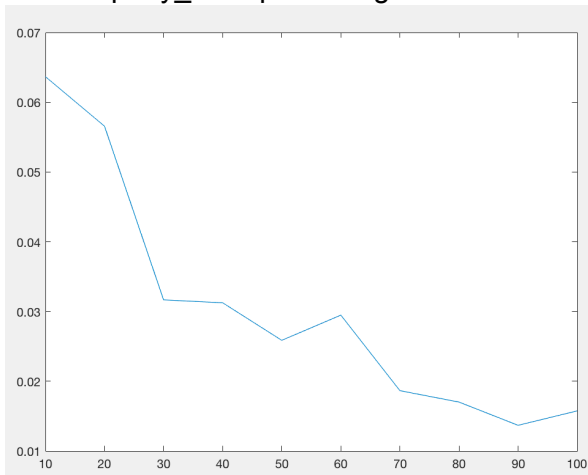
<Correct: The decimal number on above X-axis will be like below graph: 10 20 30 40 50 60 70 80 90 100>

Right Graph: y\_axis: performance

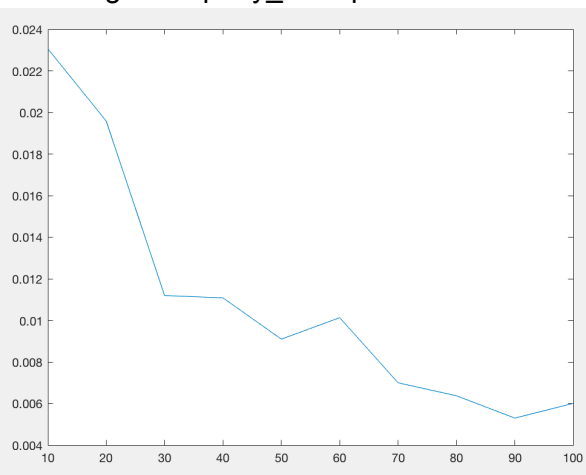


Training Function: scg

Left Graph: y\_axis: percentage Error

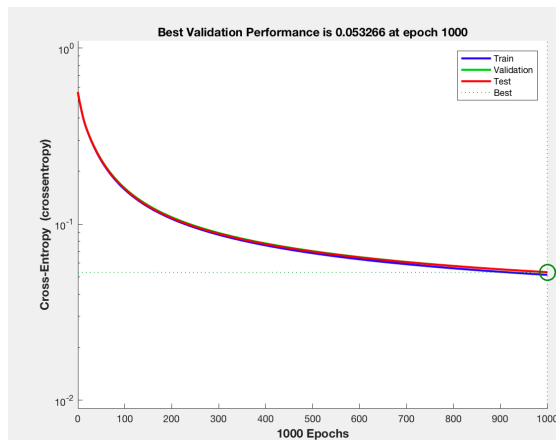


Right Graph: y\_axis: performance

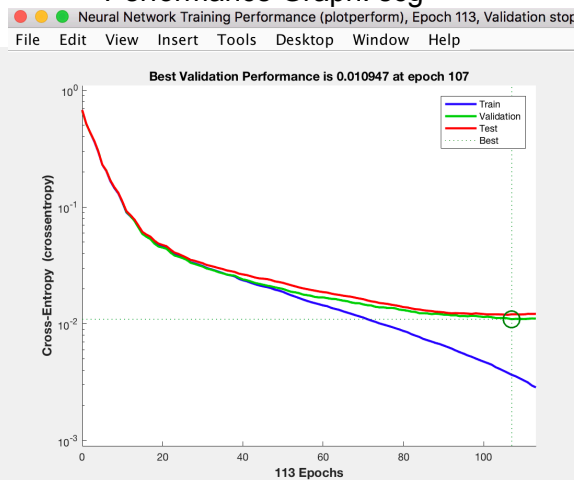


Both of **scg** and **gd** got local minimum error and performance value points around 80 to 90. Thus, in this test samples, the hidden layer set around 80 for **dg**, and around 90 for **scg** will be best.

Performance Graph: gd



Performance Graph: scg



**gd** performance during running each hidden layer shows: there is no overfitting or under fitting. And most of the performance graphs are similar.

**scg** performance during running each hidden layer shows: got overfitting in the end, when validation check reaches to 6.

### Set hidden layer size to 200(apparent larger):

**Training Function: gd**

percentage Error = 0.1389    performance = 0.0461

**Training Function: scg**

percentage Error = 0.0106    performance = 0.0043

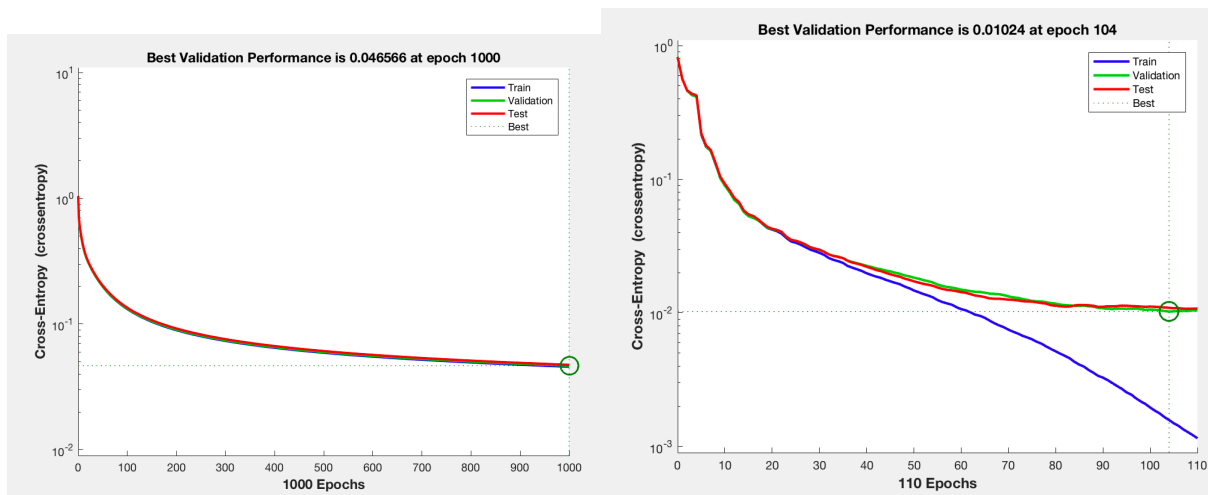
The changing of errors of gd is smoother. The changing of performance becomes not too apparent for both training functions when layer increases.

After comparing the errors and performance value of **gd** to the hidden layer which is 80 on above graph, the values are smaller, thus larger hidden layer size still affects the predict result. By now, the test sample shows 200 hidden layers could be the best choice for **gd** training function with such input data.

After comparing the errors and performance value of **scg** to the hidden layer which is 90 on above graph, the value are also about the same. Thus, 90 could be the best choice for this training function with such input data. And this error is the most lowest one for using the neural network by just using more than 90 hidden layers.

Performance Graph: gd

Performance Graph: scg



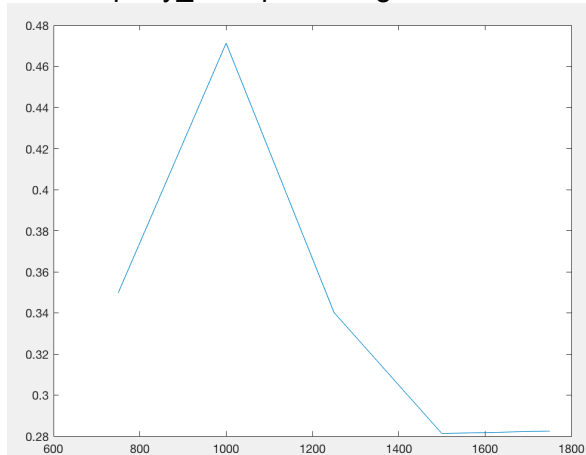
For **gd**, the performance graph is still not overfitting or undercutting.

For **scg**, the performance graph shows the overfitting in the end.

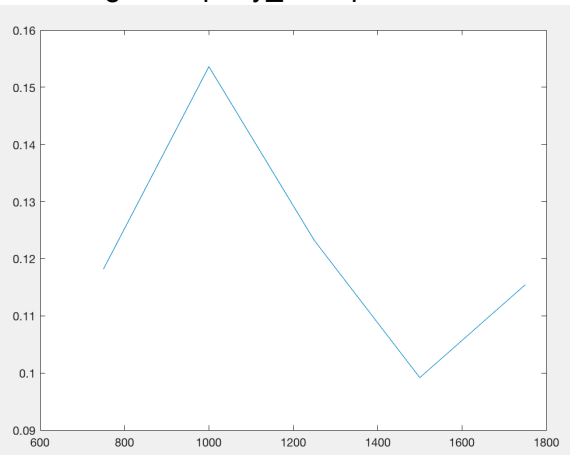
## Epochs Range(x\_axis): 750 1000 1250 1500 1750

Training Function: **gd**

Left Graph: y\_axis: percentage Error

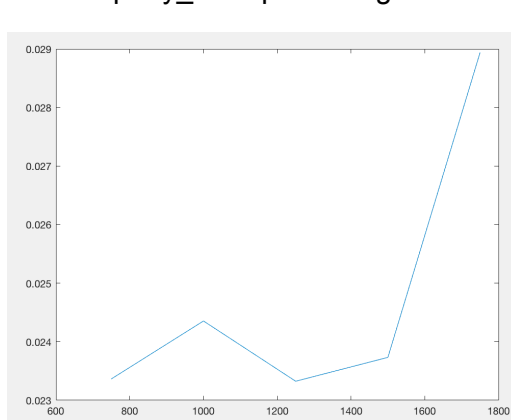


Right Graph: y\_axis: performance

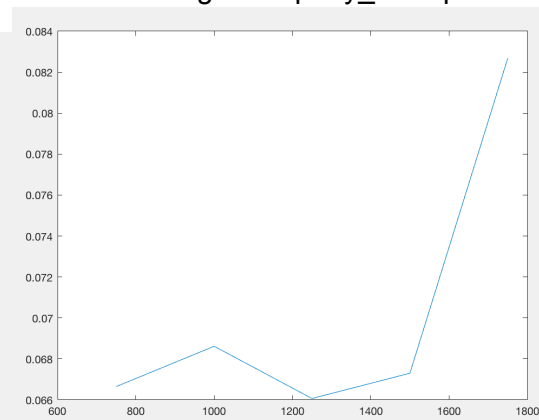


Training Function: **scg**

Left Graph: y\_axis: percentage Error



Right Graph: y\_axis: performance



For **gd**, the local minimum of error is at 1500, for **scg**, the lowest error is at 1250. Both training function shows the error becomes local maximum at epochs is 1750.

However, increasing hidden layer size still improves the neural performance better than changing the other two features.