

BT3041 Analysis and Interpretation of Biological Data

BT3041 AIBD Assignment 1

N Sowmya Manojna | BE17B007
Department of Biotechnology,
Indian Institute of Technology, Madras



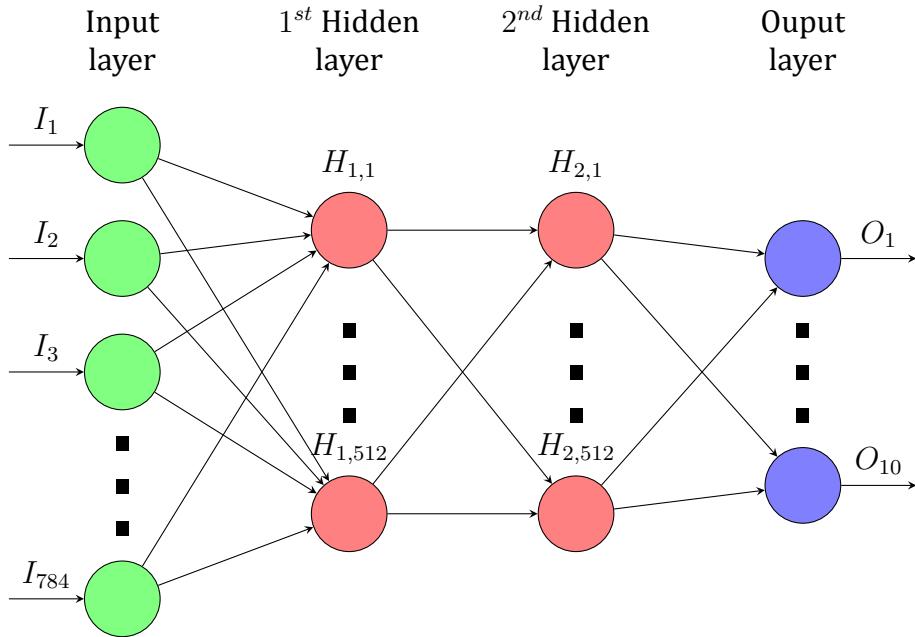
1 Question 1 - MLP

The solution to this question is coded using MATLAB.

Codes for this question have been written using MATLAB 2020a which performs automatic remapping of matrices or vectors in matrix operations.

1.1 Network Architecture

The neural architecture is as follows:



The network comprises of:

- A input layer with 784 neurons
- Two hidden layers comprising of 512 neurons each
- An output layer with 10 neurons

1.2 Activation functions

- All the activation functions used are logistic activation functions.
- The sigmoid midpoint is placed at 0 and
- The logistic growth rate is made 15 fold.

The function used is:

$$f(x; c1, c2) = \frac{1}{1 + \exp^{-c1(x-c2)}}$$

Where, $c1$ denotes the midpoint and $c2$ denotes the growth rate of the function.

A steep growth rate is used because, most of the $z1, z2$ values lie between the range -1 and 1. The sigmoidal function is largely linear in this phase and prevented the network from learning non-linear relations. Hence, in order to introduce non-linearity the logistic growth rate was increased and the parameters used were $c1=15$, $c2=0$.

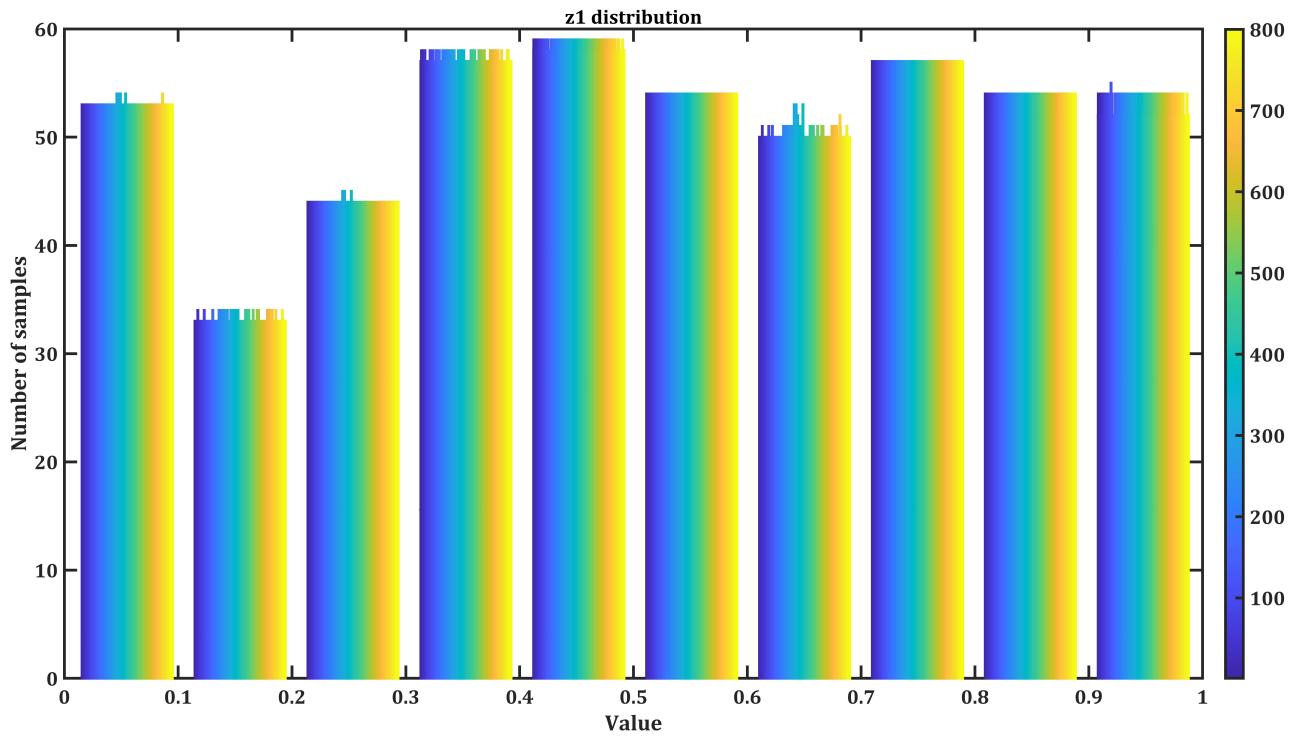


Figure 1: Distribution of values in $z1$

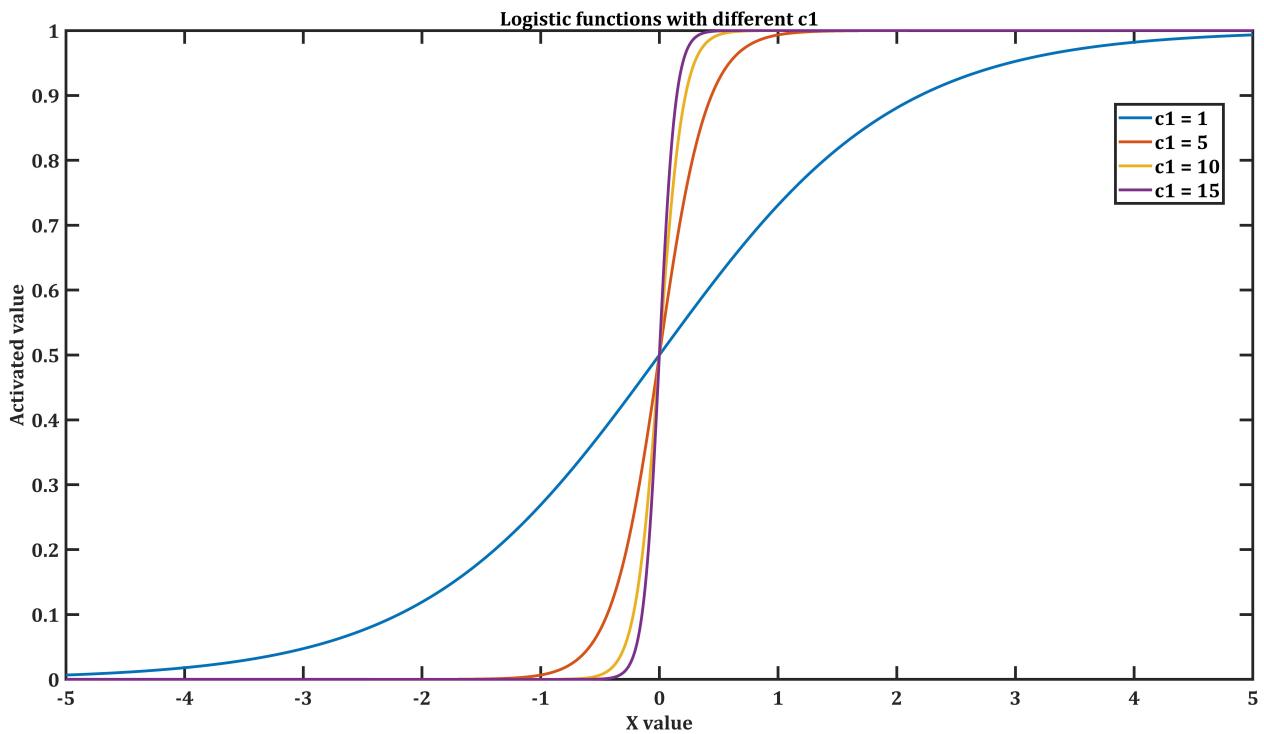


Figure 2: Plots of logistic functions with varying slope

The activation function with the constants is:

$$f(x) = \frac{1}{1 + \exp^{-15x}}$$

1.3 Loss Function

Multi-class Mean Squared Error has been used as the loss function.

The equation for the same is (vector notation):

$$\text{loss} = \frac{1}{2} \|(\text{predicted_class_vector} - \text{true_class_vector})\|^2$$

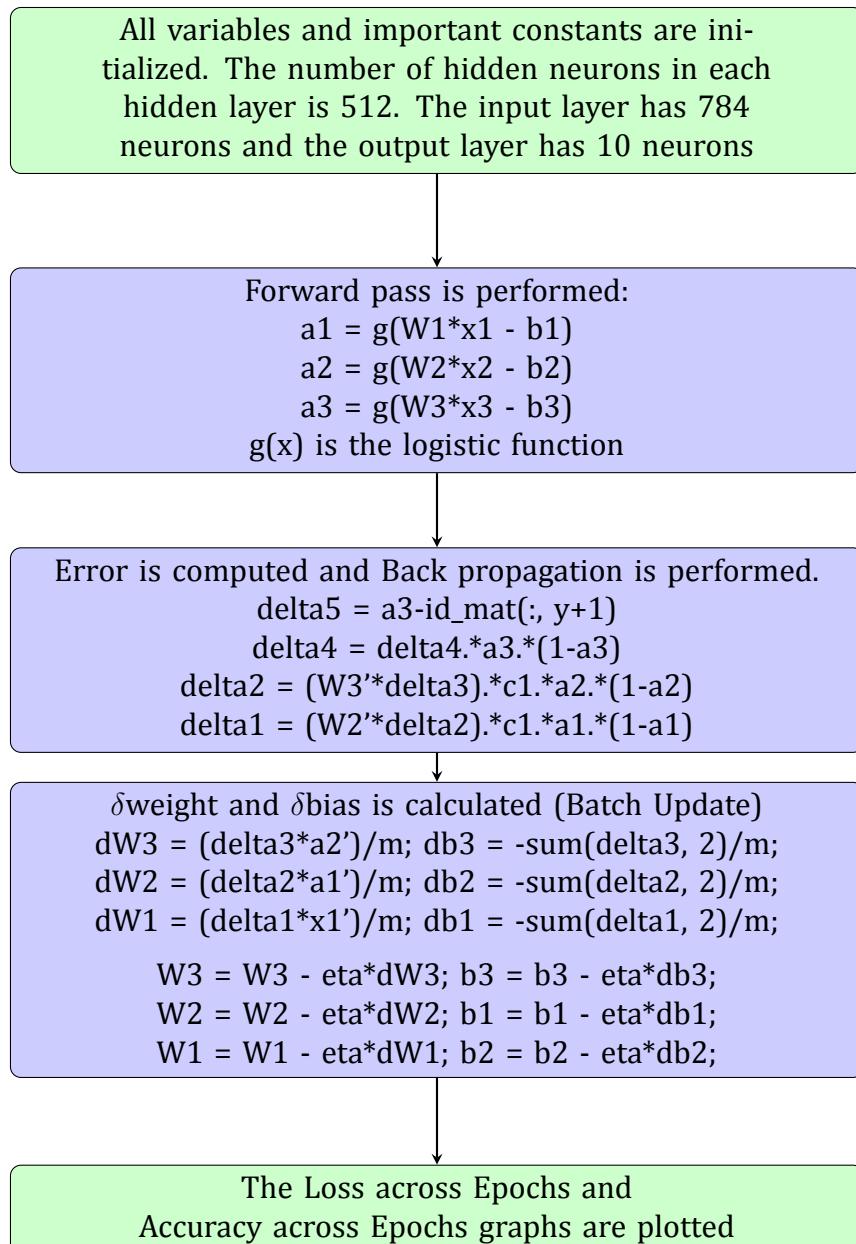
The `predicted_class_vector` and the `true_class_vector` are both of size $10 * 1$ where, an entry having 1 indicates that the predicted or actual number (i.e.)

- $(1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0)^T$ means 0
- $(0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0)^T$ means 1
- $(0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0)^T$ means 2 and so on

1.4 Learning Algorithm

The multi-class feed-forward network uses back-propagation for learning.

The flowchart of the algorithm is as follows:



The detailed algorithms is as follows:

```
1 % Accessing the images, labels and normalizing the images
2 x1 = double(trainImages')/255;
3 y = double(trainLabels');
4 t1 = double(testImages')/255;
5 yt = double(testLabels');
6
7 % Initializing the number of train and test samples
8 m = 800;
9 mt = 200;
10
11 % Defining all essential constants
12 % - eta_max and eta_min for the exponential eta decay function
13 % - epochs for training
14 % - fact (factor) for determinig the rate of eta decay
15 % - c1 rate of logistic function
16 % - save_every to save the accracy values after a fixed number of epochs
17 eta_max = 0.3;
18 eta_min = 0.01;
19 epochs = 6000;
20 fact = 5000;
21 c1 = 15;
22 c2 = 0;
23 save_every = 10;
24
25 % Defining the size of the network
26 n1 = size(x1,1);
27 n2 = 512;
28 n3 = 512;
29 n4 = 10;
30
31 % Initializing the weights, biases and additional variables
32 % to store the accuracy and
33 W1 = zeros(n2, n1);
34 W2 = zeros(n3, n2);
35 W3 = zeros(n4, n3);
36
37 b1 = rand(n2, 1);
38 b2 = rand(n3, 1);
39 b3 = rand(n4, 1);
40
41 id_mat = eye(n4);
42 train_accuracy = zeros(ceil(epochs/save_every),1);
43 test_accuracy = zeros(ceil(epochs/save_every),1);
44
45 % Start training
46 for i = 1:epochs
47     % Forward pass
48     z1 = W1*x1 - b1;
49     a1 = sigmf(z1, [c1, 0]);
50
51     z2 = W2*a1 - b2;
52     a2 = sigmf(z2, [c1, 0]);
53
54     z3 = W3*a2 - b3;
55     a3 = sigmf(z3, [1, 0]);
56
57     % Train and test accuracy calculation
58     if rem(i, save_every) == 0
```

```

59     fprintf('Epoch: %d; eta: %f \n', i, eta);
60     [~, idx] = max(a3);
61     train_accuracy(i/save_every) = sum(idx' == y+1)/m;
62
63     % Testing
64     t2 = sigmf(W1*t1 - b1, [c1,0]);
65     t3 = sigmf(W2*t2 - b2, [c1,0]);
66     t4 = sigmf(W3*t3 - b3, [1,0]);
67     [~, idx] = max(t4);
68     test_accuracy(i/save_every) = sum(idx' == yt + 1)/mt;
69 end
70
71 % Backpropagation
72 d4 = a3 - id_mat(:, y+1);
73 d3 = d4.*c1.*a3.*(1-a3);
74 d2 = (W3'*d3).*c1.*a2.*(1-a2);
75 d1 = (W2'*d2).*c1.*a1.*(1-a1);
76
77 dW3 = (d3*a2')/m;
78 dW2 = (d2*a1')/m;
79 dW1 = (d1*x1')/m;
80 db3 = -sum(d3, 2)/m;
81 db2 = -sum(d2, 2)/m;
82 db1 = -sum(d1, 2)/m;
83
84 % get eta based on epoch
85 eta = get_exp_eta(i, eta_max, eta_min, fact);
86
87 % Weight and bias update
88 W3 = W3 - eta*dW3;
89 W2 = W2 - eta*dW2;
90 W1 = W1 - eta*dW1;
91 b3 = b3 - eta*db3;
92 b2 = b2 - eta*db2;
93 b1 = b1 - eta*db1;
94 end

```

1.5 Learning Rate

A combination of exponential decay and step learning rate schedule is used. The initial learning rate (maximum learning) rate is 0.2 and the final learning rate (minimum) is 0.01.

The learning rate schedule is plotted below:

The learning rate schude is as follows:

$$\eta = \begin{cases} \eta_{min} + (\eta_{max} - \eta_{min}) * \exp^{\frac{-epochs}{fact*0.5}} & epochs < fact \\ \eta_{min} & otherwise. \end{cases}$$

The η_{min} and η_{max} used in this model are 0.01 and 0.3 respectively.

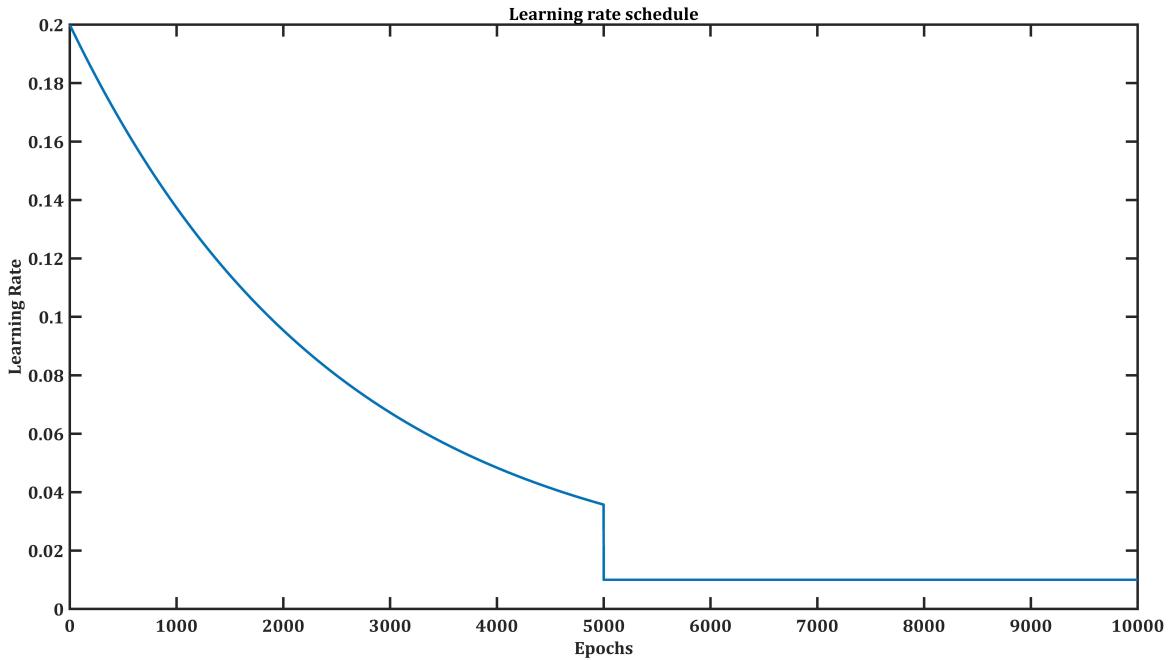


Figure 3: Learning rate schedule used in the model

1.6 Loss and Accuracy testing

The accuracy obtained after training for 6000 epochs is as follows:

- Train accuracy: 99.5%
- Test accuracy: 87%

The final loss obtained after training for 6000 epochs is 0.0108. The loss distribution across epochs is as follows:

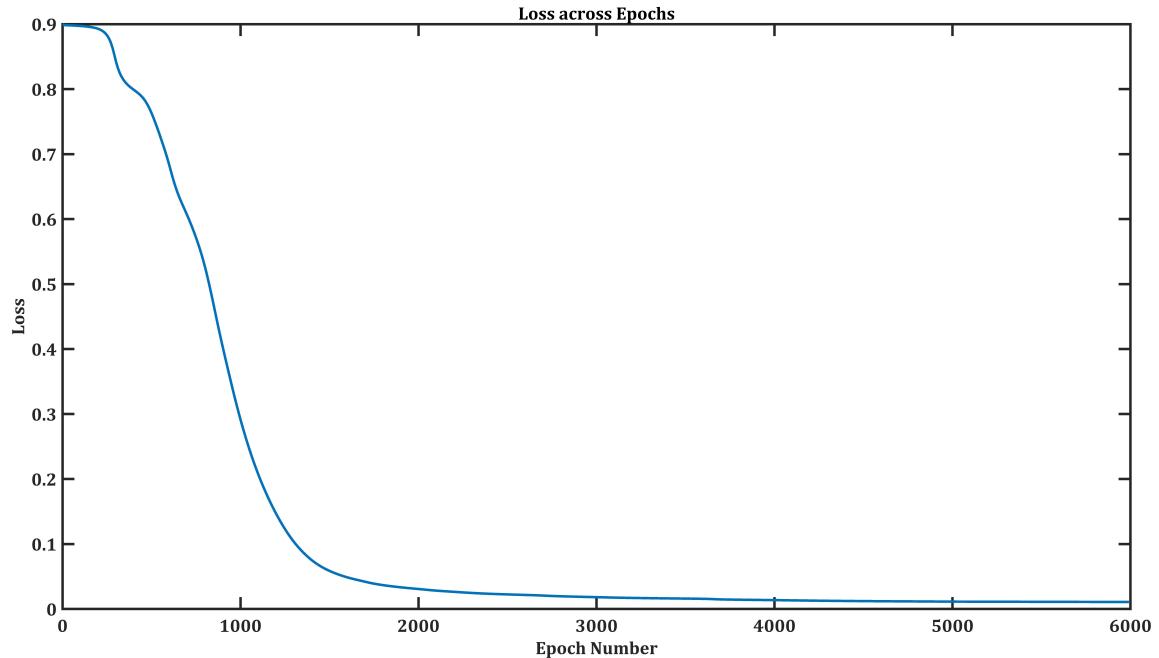


Figure 4: Loss function value across epochs

The test and train accuracy distribution across epochs is as follows:

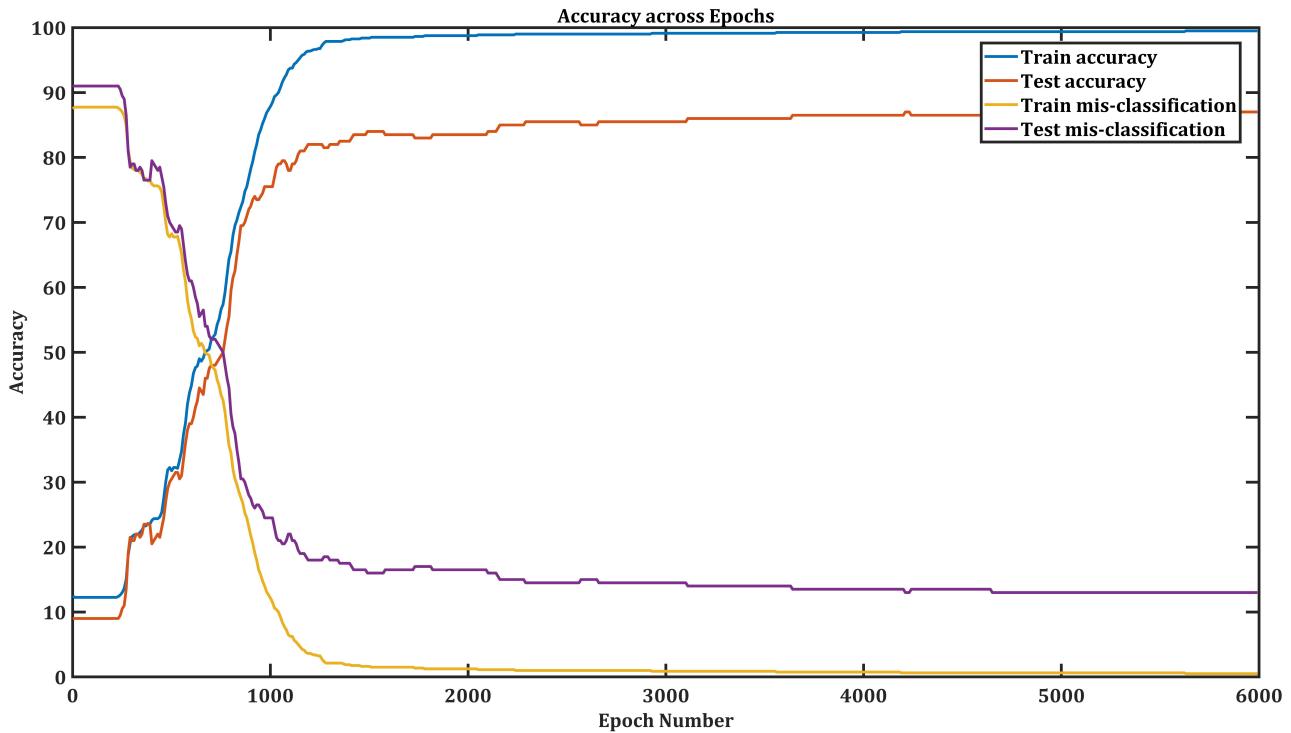


Figure 5: Accuracy across epochs

2 Question 2 - DBSCAN

The solution to this problem is coded using Python3 and the Jupyter framework is used for the same.

All the libraries used for this question can be installed by using the following code (inside question_2 folder): `$ pip3 install -r requirements.txt`

2.1 Parameters

The parameters used to obtain 4 clusters is:

- Epsilon radius (ϵ) = 0.3
- `min_points` = 40

2.2 Procedure

- The distance between all pairs of points was calculated.
- The number of neighbors a point has is calculated by considering all distances that are less than `eps`.
- Points that have more than `min_points` number of neighbors are labeled **core points**.
- Points that have less than `min_points` number of neighbors, but are reachable from core points are labeled **neighbor points**.
- Points that have less than `min_points` number of neighbors and are not neighbor points are labeled **noise points**.

- Cluster allocation is done using a Breadth First Search (BFS) like spanning through the core points. All the accessible core points and their neighbors are assigned the same cluster number.
- Once all the accessible core points and their neighbors are given a cluster label, the current cluster label is incremented and the cluster allocation is done using the next set of core points.

2.3 Results

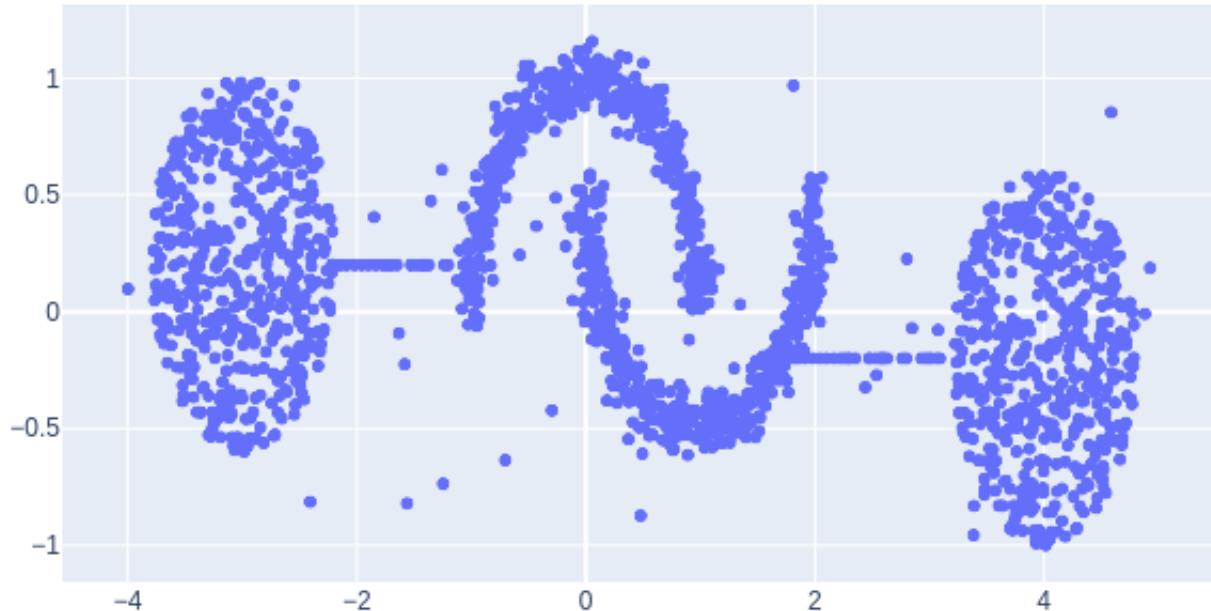


Figure 6: The initial unclustered data

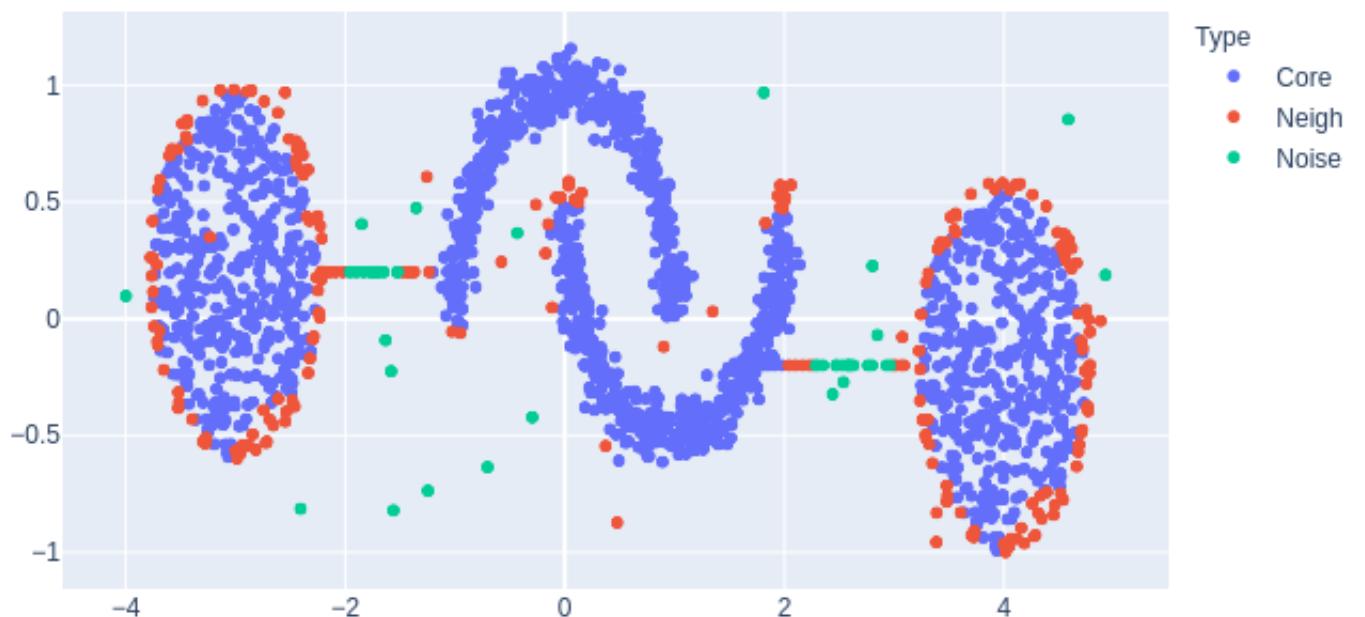


Figure 7: Core, Neighbor and Noise points labeled

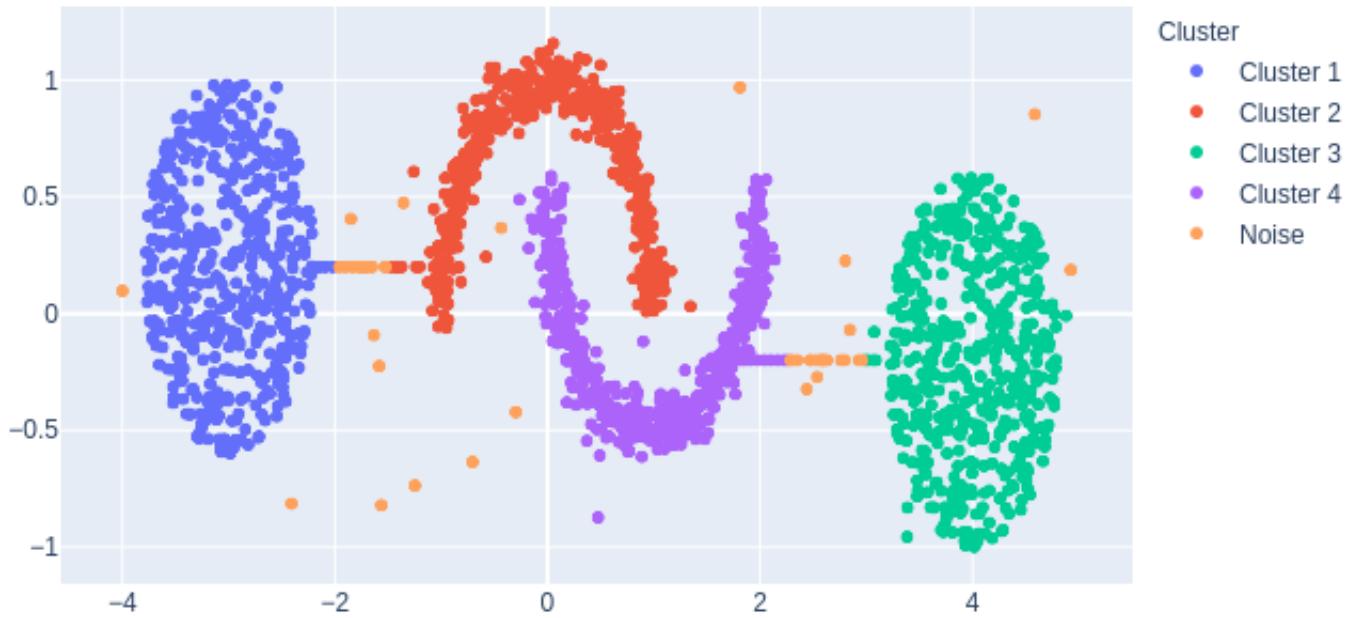


Figure 8: Final clustered data

The distribution of points in clusters is as follows:

- Cluster 1: 509 points
- Cluster 2: 468 points
- Cluster 3: 512 points
- Cluster 4: 466 points
- Noise: 45 points