

Pavan Rathnakar Shetty

EECE 5644

HW 4

Due date: December 14, 11:59 pm

GitHub: Files in the hw4 folder

<https://github.com/Pavan-r-shetty/5644.git>

## Question 1

### Problem 1

Multi-layer Perceptron was used to approximate class label posteriors.

Loss used for training < Minimum average cross-entropy loss.

Trained models were then used to approximate a MAP classification rule to achieve minimum probability of error on validation dataset.

Data: 4 classes with uniform priors

class 0:

$$p(L=0) = 0.25, \quad m_0 = \begin{bmatrix} 4 \\ 2 \\ 1 \end{bmatrix} \quad \text{cov } 0 = \begin{bmatrix} 1 & 0 & 0.2 \\ 0 & 0.1 & 0.3 \\ 0.5 & 0.6 & 0.8 \end{bmatrix}$$

class 1:

$$p(L=1) = 0.25, \quad m_1 = \begin{bmatrix} 1 \\ 3 \\ 4 \end{bmatrix} \quad \text{cov } 1 = \begin{bmatrix} 1 & 0 & 0.6 \\ 0.6 & 1 & 0 \\ 0 & 0.6 & 0.8 \end{bmatrix}$$

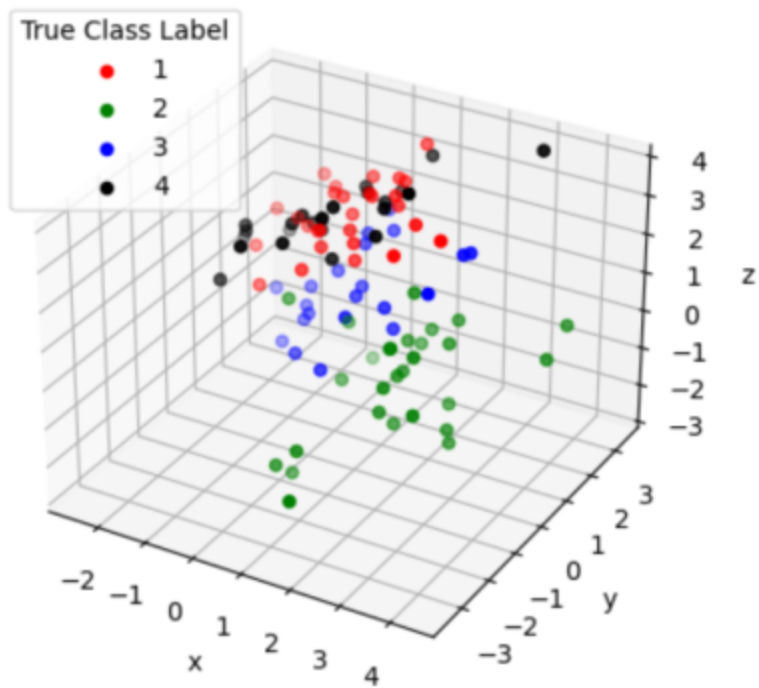
class 2:

$$p(L=2) = 0.25, \quad m_2 = \begin{bmatrix} 2 \\ 4 \\ 2 \end{bmatrix} \quad \text{cov } 2 = \begin{bmatrix} 0.8 & 0.1 & 0.7 \\ 0.4 & 0.9 & 0.6 \\ 0.7 & 0.8 & 0.9 \end{bmatrix}$$

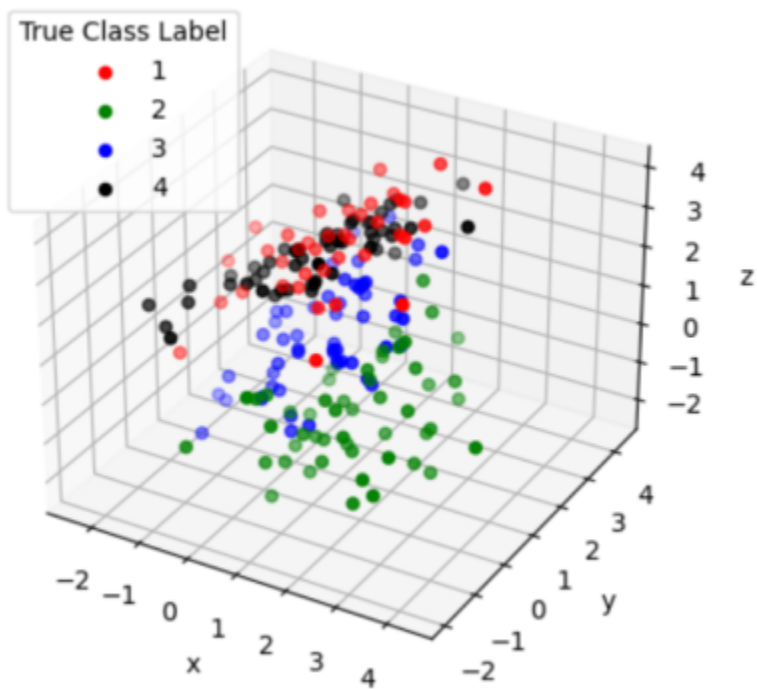
class 3:

$$p(L=3) = 0.25, \quad m_3 = \begin{bmatrix} 3 \\ 1 \\ 3 \end{bmatrix} \quad \text{cov } 3 = \begin{bmatrix} 1 & 0.6 & 0.9 \\ 0.1 & 0.8 & 0.3 \\ 0.7 & 0.3 & 1 \end{bmatrix}$$

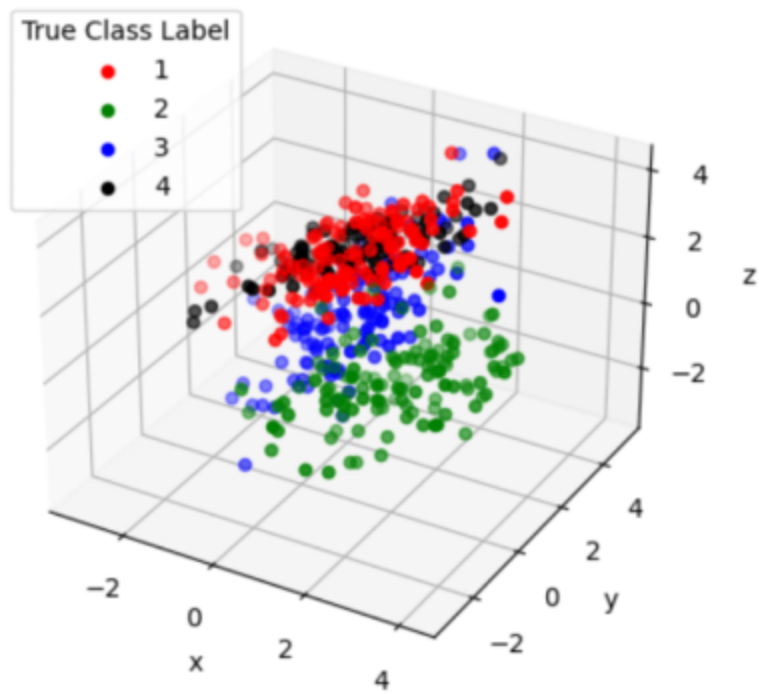
The scattered plots for all the dataset are?



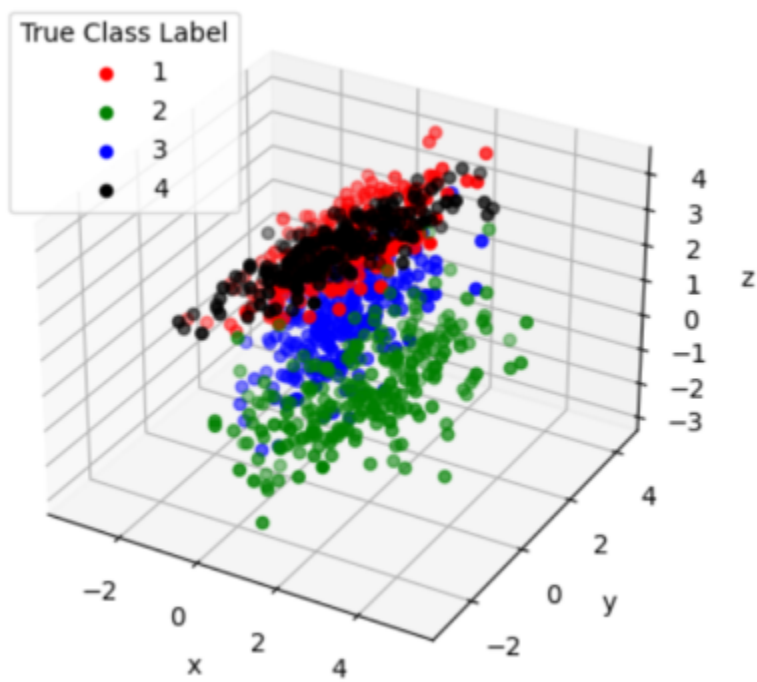
train\_data\_100



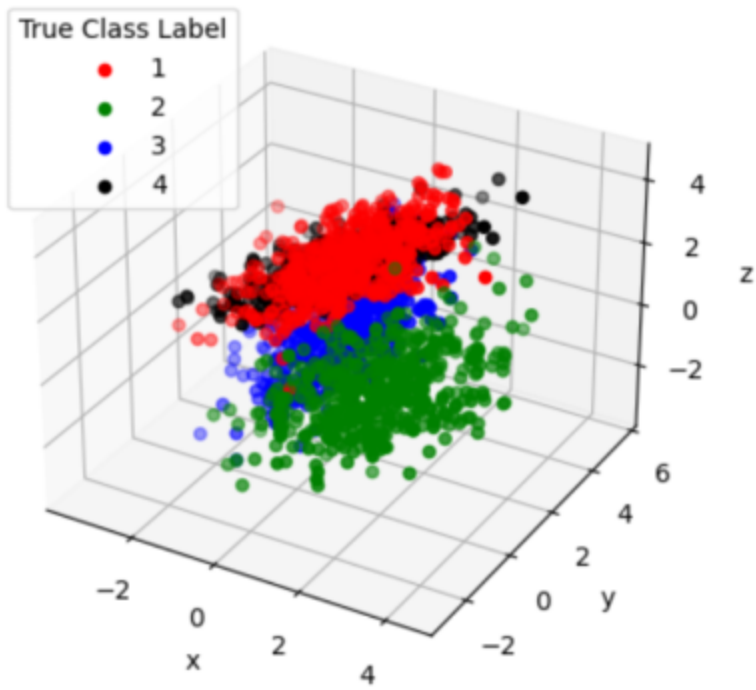
train\_data\_200



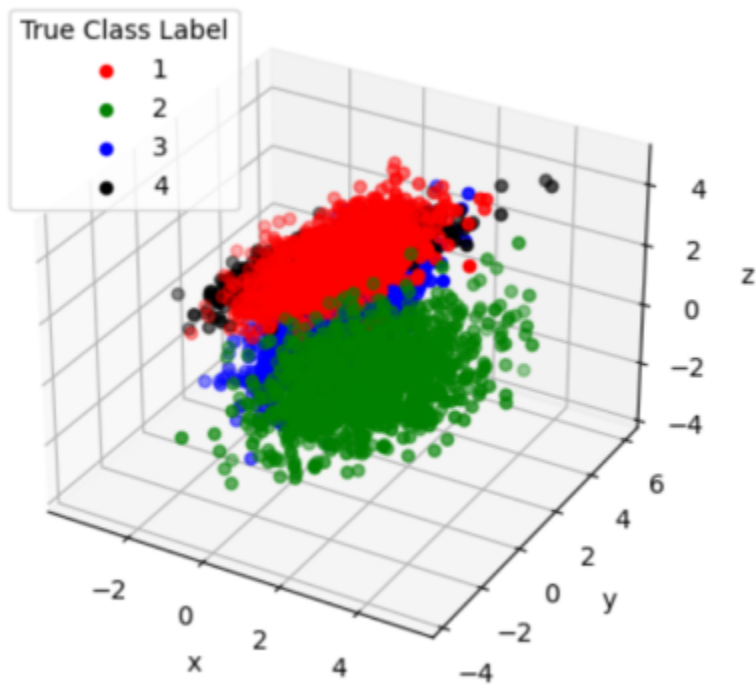
train\_data\_500



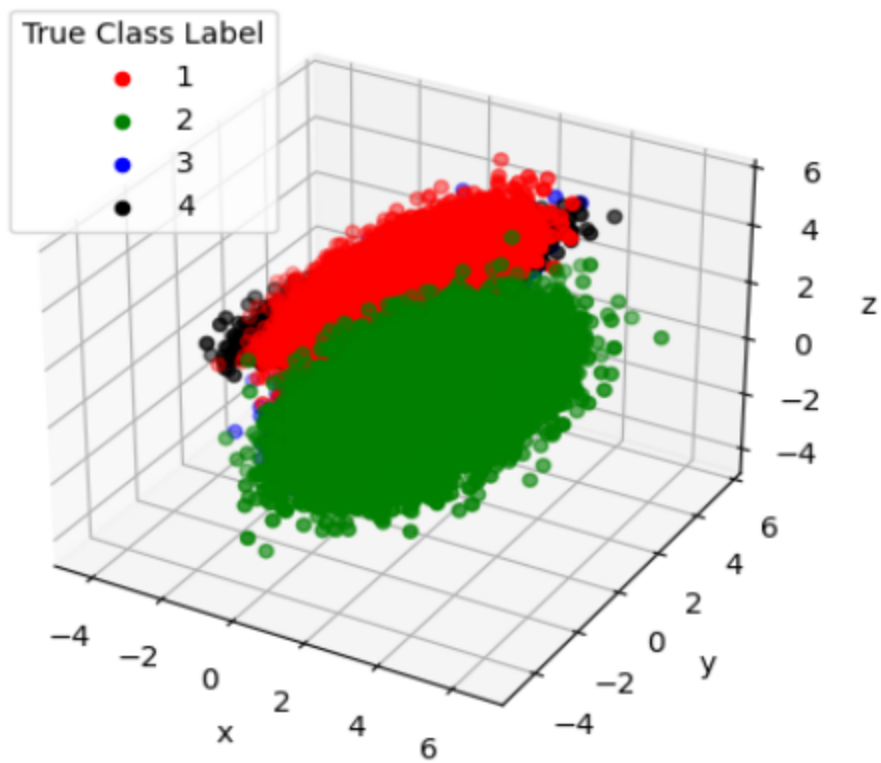
train\_data\_1000



train\_data\_2000



train\_data\_5000



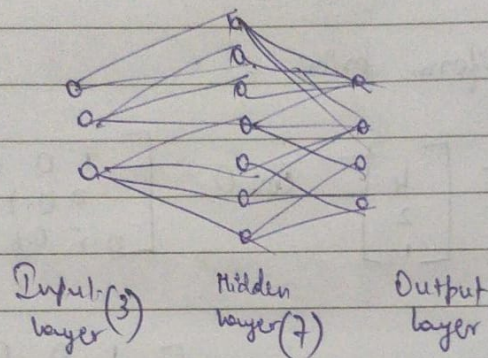
train\_data\_10000



## Multi layer Perceptron:

- A model is created with an input layer, one hidden layer and output layer.
- The number of nodes in the hidden layer is varied from 1 to 10 to select the best.

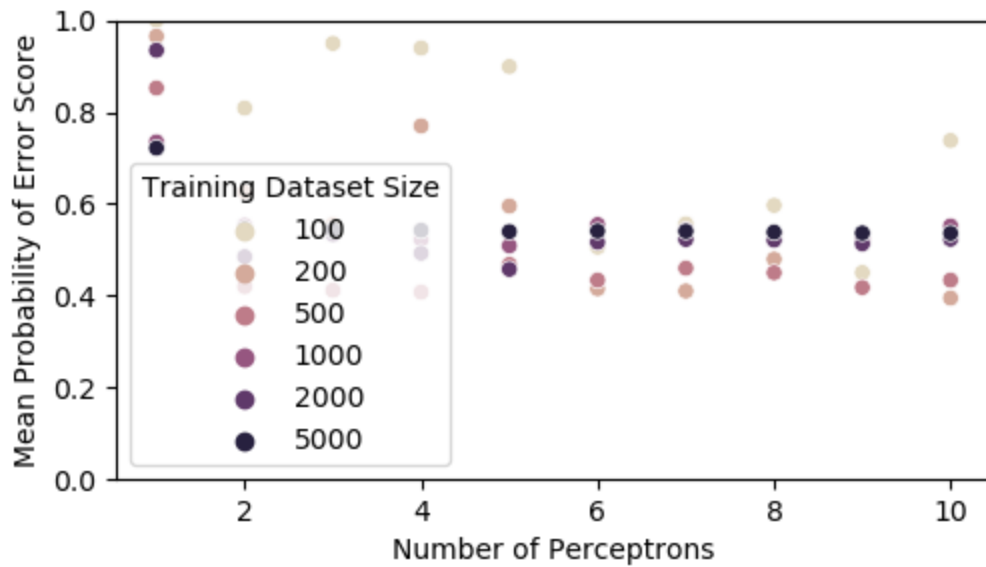
Eg:



## Hyper parameter tuning using 10-fold cross validation

Figure below shows the performance of MLP with varying number of perceptrons.

- 10 fold cross validation was performed to determine the optimal number of MLP.



### Observation and Inference:

- As perceptron increases, error reduced, however for 4-6 perceptron there is no much significant reduction in the error score.
- The score saturated at sometimes overfit.
- Larger the training set, better is the performance.

- Optimal Perceptron (minimum error) was trained on all the dataset (100 - 500 train data)

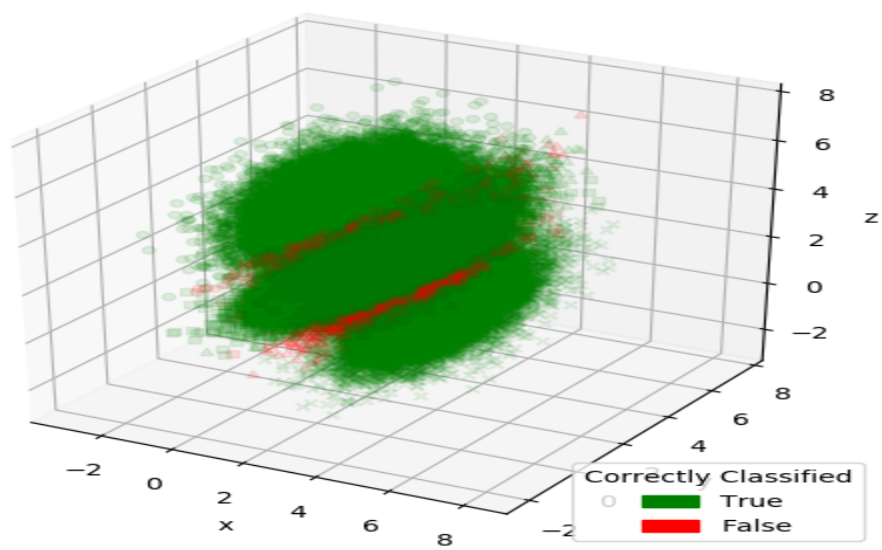
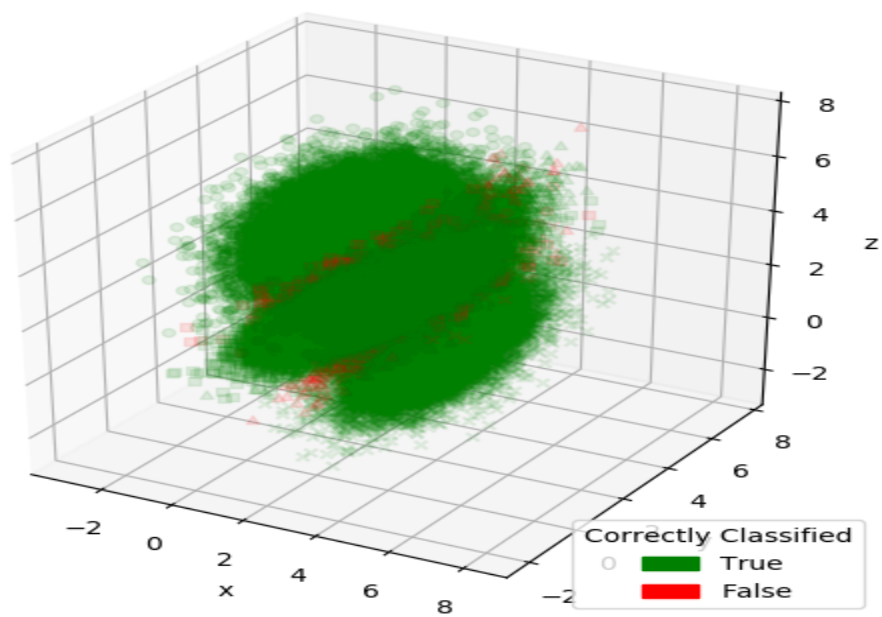
- Theoretical optimal classifier was used using minimum probability of error classification rule.

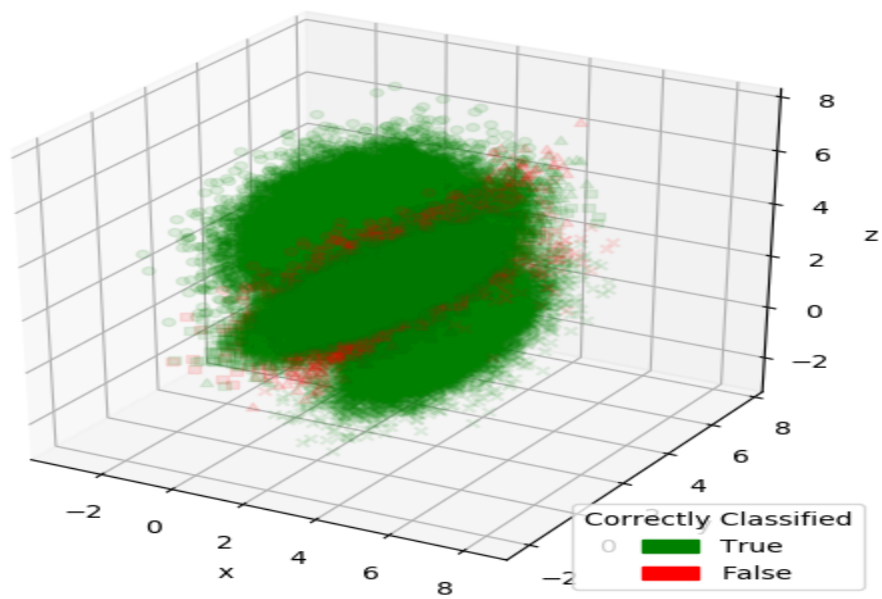
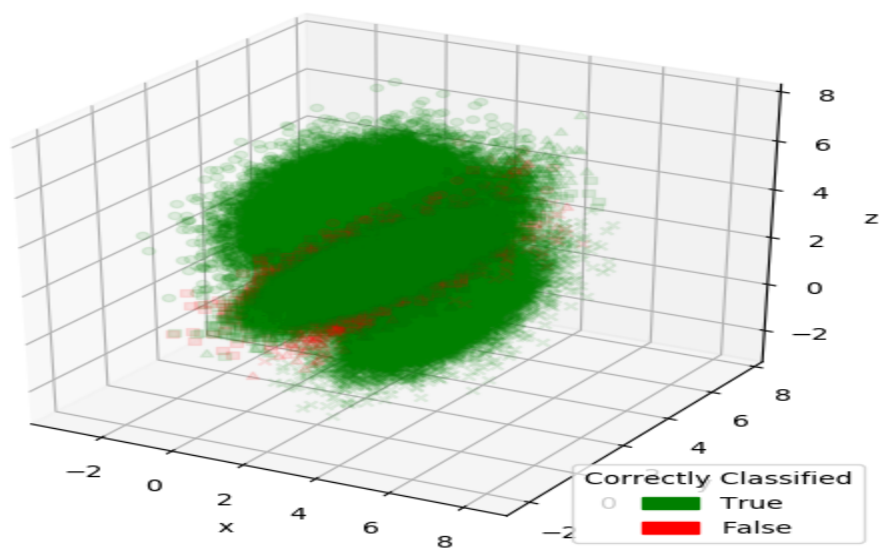
- Theoretical optimal error = 16.8%.

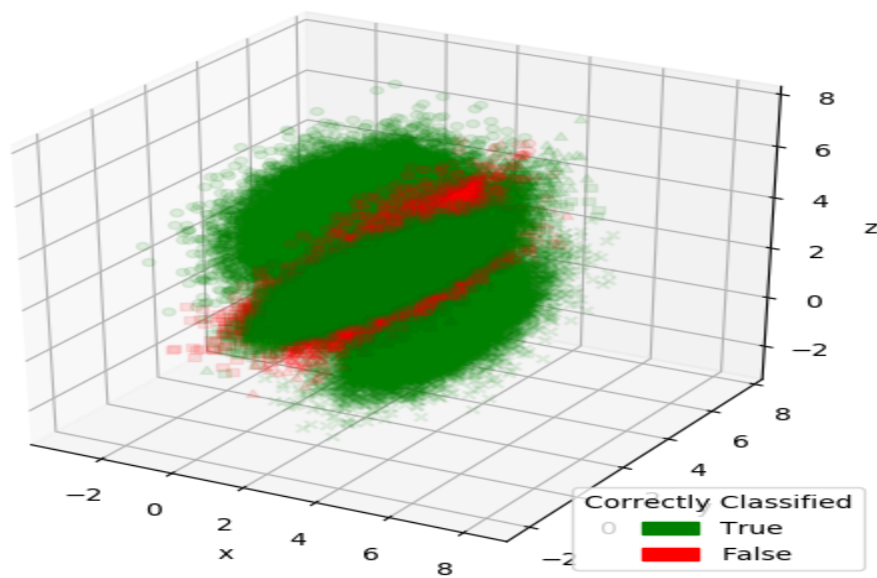
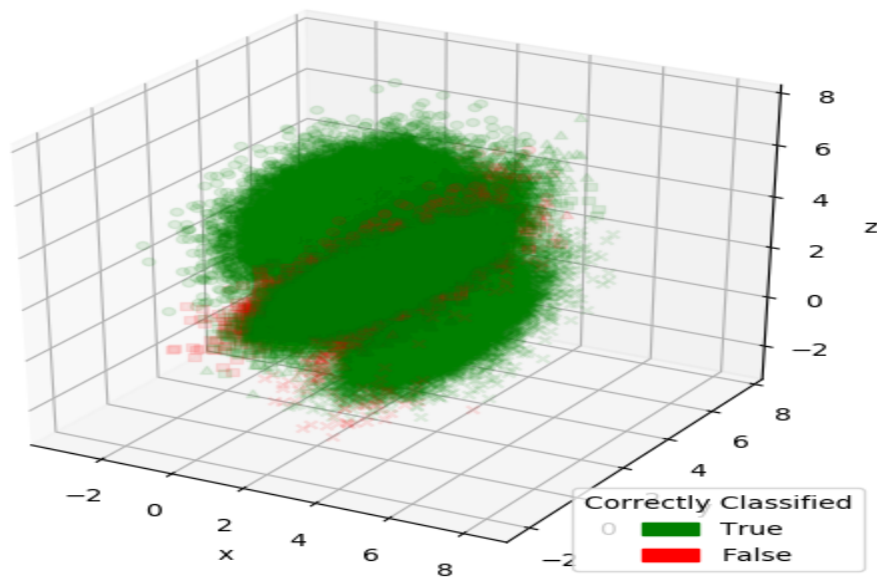
Figure below shows the classification done by the model on datasets. data-10,000 was used to validate.

Classification with 5000-100 Train data is shown below

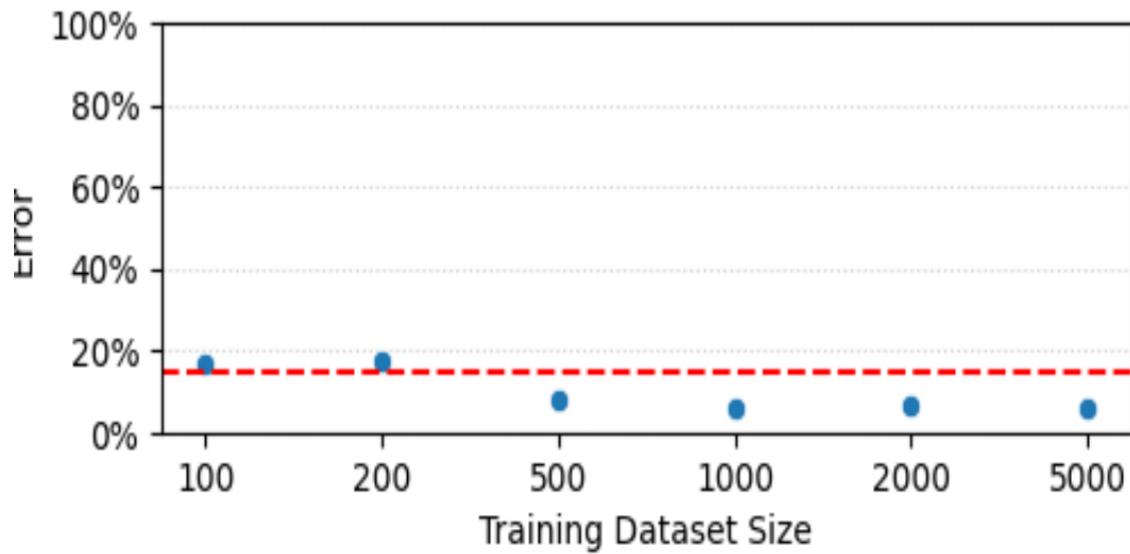








The figure below shows error vs training data size for each MAP along with MAP classifier. (Red line = 16.8% optimal classifier)



#### Summary of Performance

Model	Size	Probability of Error (%)
train_data_100	100	17.12
train_data_200	200	0.175
train_data_500	500	7.5
train_data_1000	1000	6.4
train_data_2000	2000	6.7
train_data_10000	10000	6.1

Packages Used: Keras, Scikit Learn, Pandas, Matplotlib



## Question 2:

### Problem 2:

Generating Gaussian Mixture Model for 2D data should have 4 components with unique Probability of selection, means and covariances.

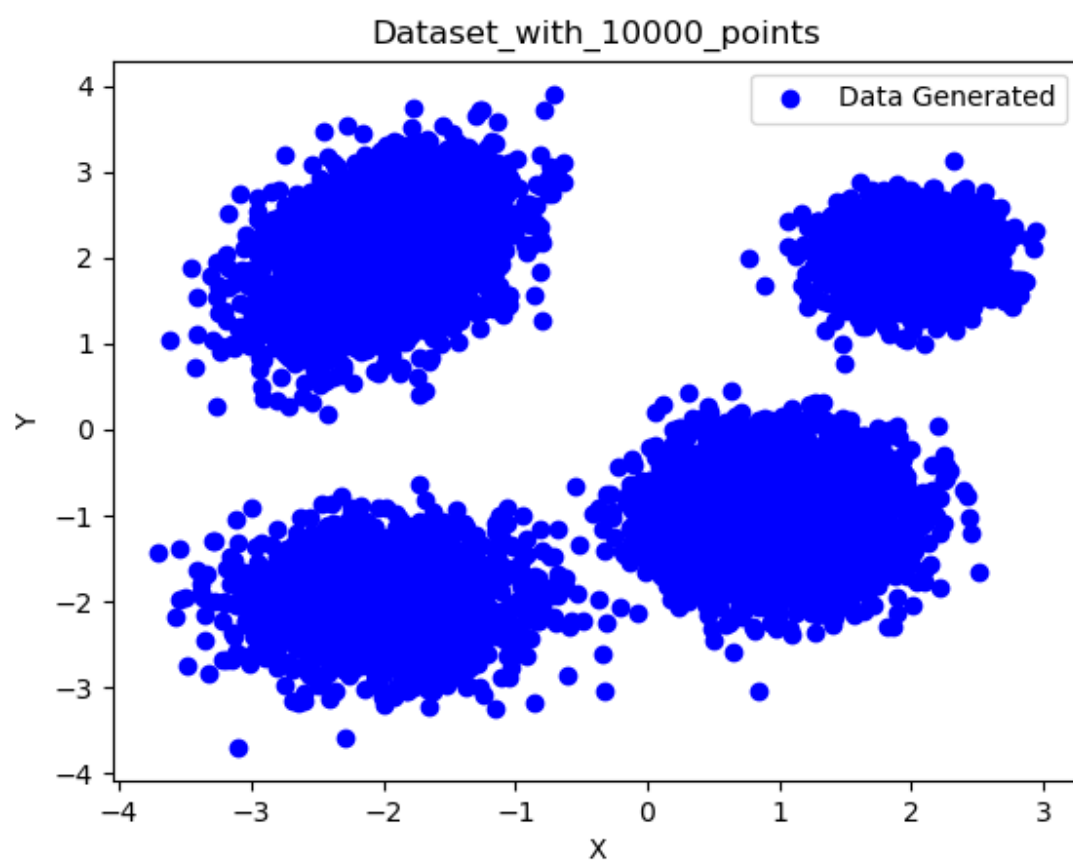
The  $\mu$ 's,  $\sigma$ 's & weight chosen are:

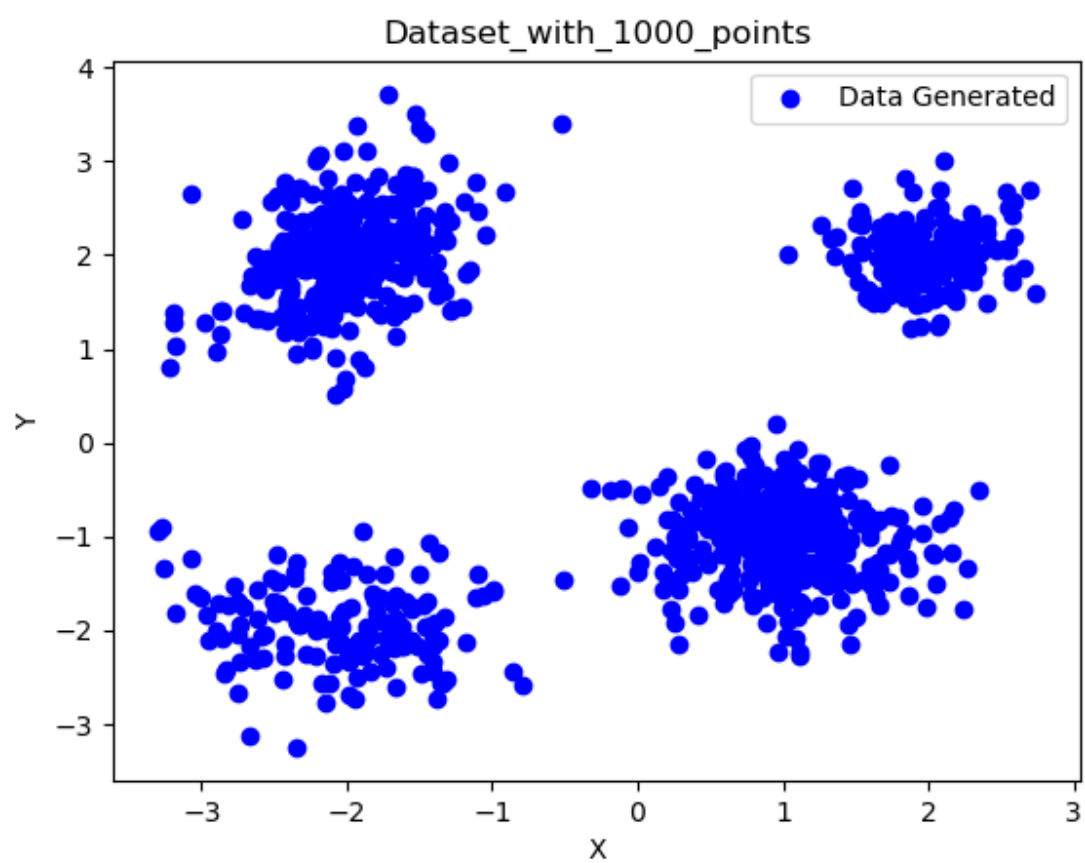
$$w_0 = 0.2, \mu_0 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \text{cov}_0 = \begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix}$$

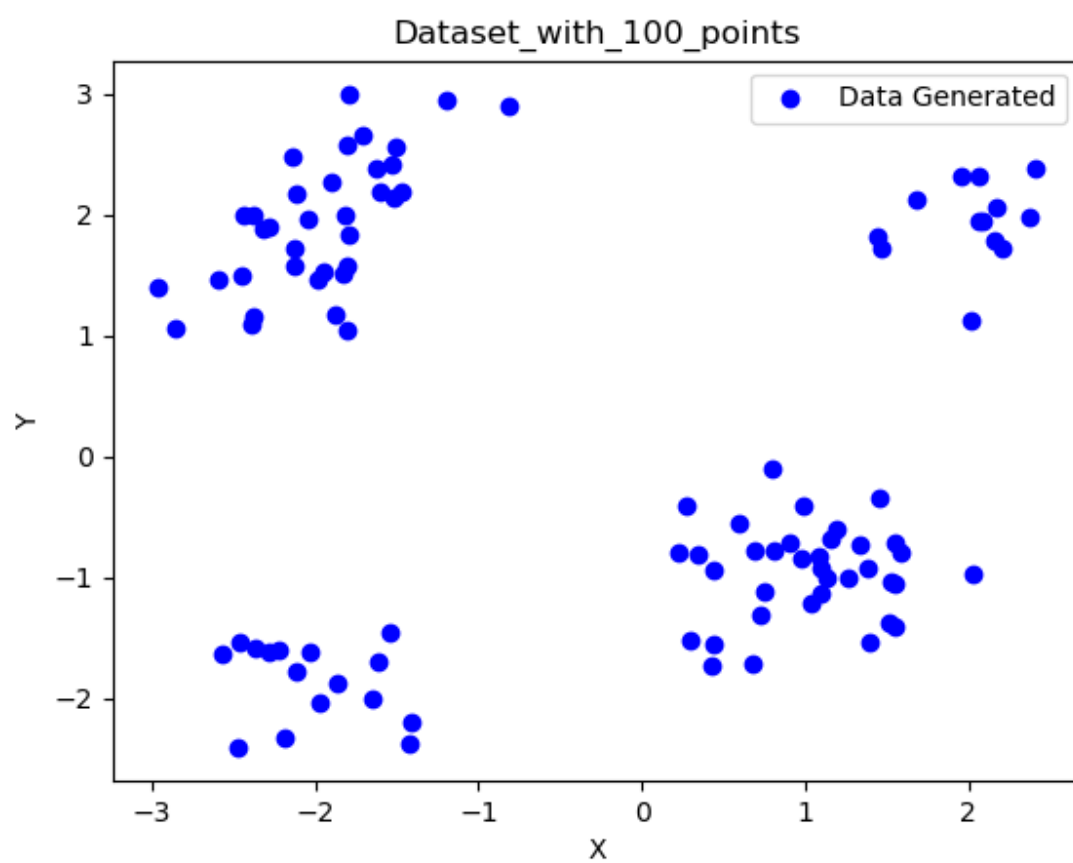
$$w_1 = 0.3, \mu_1 = \begin{bmatrix} -2 \\ 2 \end{bmatrix}, \text{cov}_1 = \begin{bmatrix} 0.2 & 0.1 \\ 0.1 & 0.3 \end{bmatrix}$$

$$w_3 = 0.15, \mu_2 = \begin{bmatrix} -2 \\ -2 \end{bmatrix}, \text{cov}_2 = \begin{bmatrix} 0.3 & 0 \\ 0 & 0.2 \end{bmatrix}$$

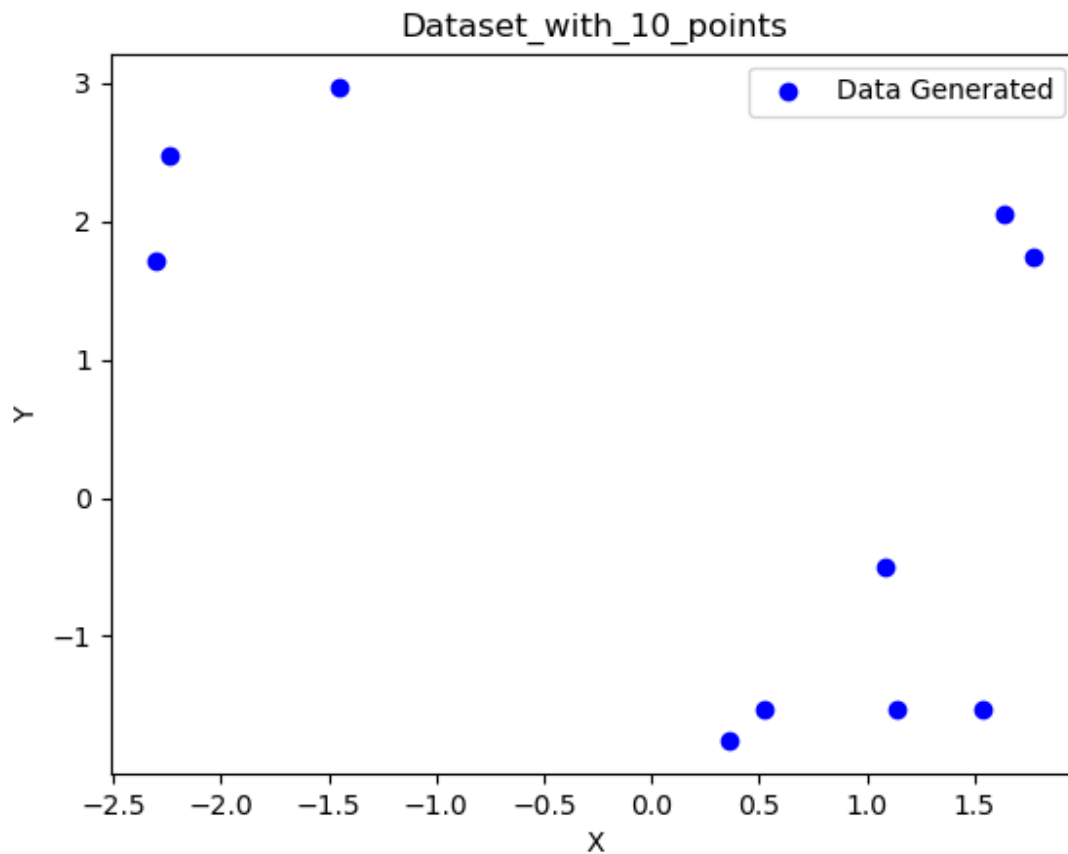
$$w_4 = 0.35, \mu_3 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \text{cov}_3 = \begin{bmatrix} 0.2 & 0 \\ 0 & 0.2 \end{bmatrix}$$











For each dataset using ML $\bar{E}$  and 10-fold cross-validation, the gaussian components were evaluated with different model ranging 1-6.

The Bayesian Information Criteria and Akaike's Information Criteria were also calculated with -ve log likelihood.

The experiment was carried out 30 different times for each dataset. The frequency was noted down.

**10 points dataset:**

<b>GMM order selection based on (-) log-likelihood</b>	<b>Frequency</b>
<b>1</b>	<b>30</b>

**100 points dataset:**

<b>GMM order selection based on (-) log-likelihood</b>	<b>Frequency</b>
<b>4</b>	<b>25</b>
<b>5</b>	<b>5</b>

**1000 points dataset:**

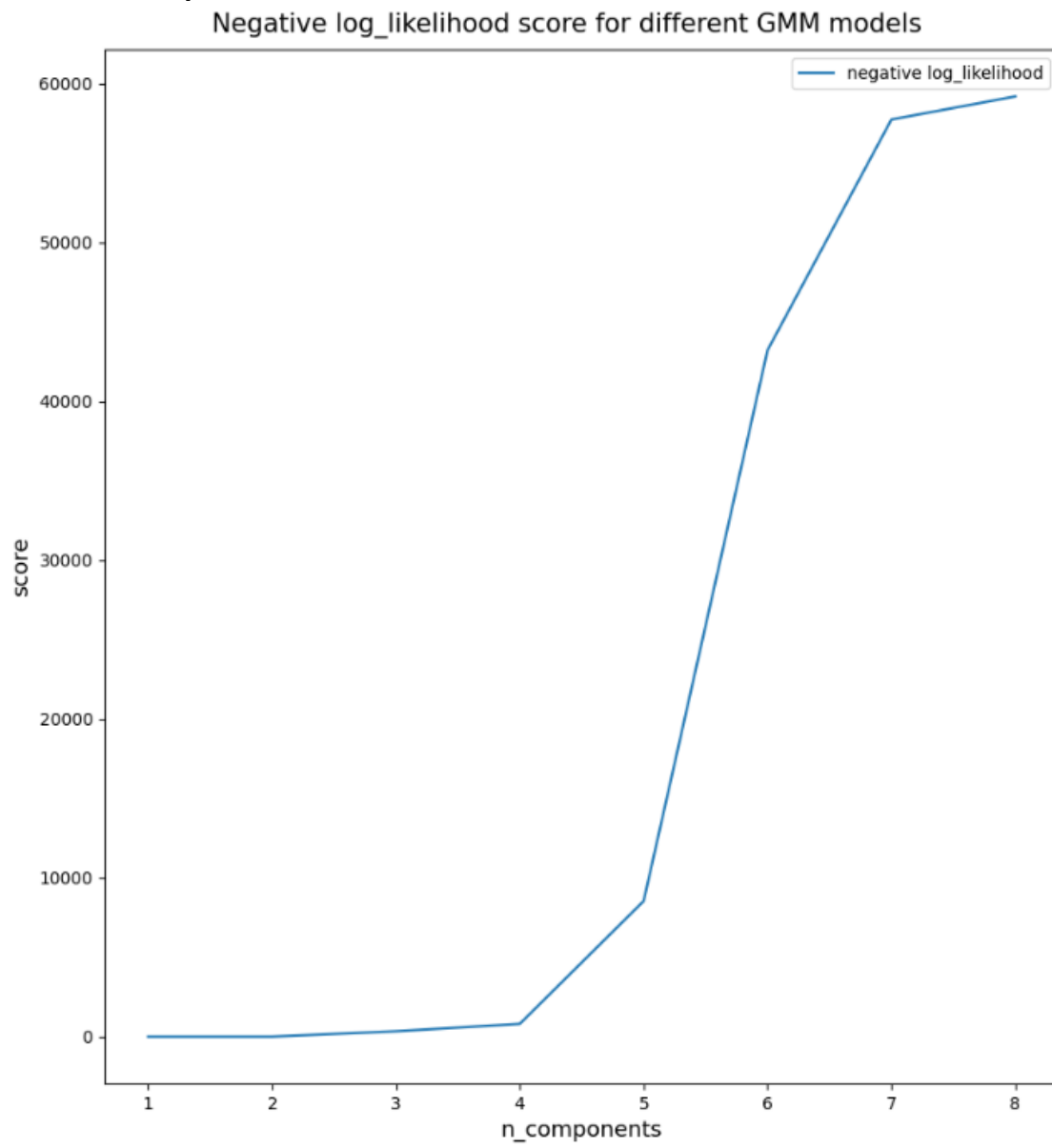
<b>GMM order selection based on (-) log-likelihood</b>	<b>Frequency</b>
<b>4</b>	<b>23</b>
<b>5</b>	<b>5</b>
<b>6</b>	<b>2</b>

**10000 points dataset:**

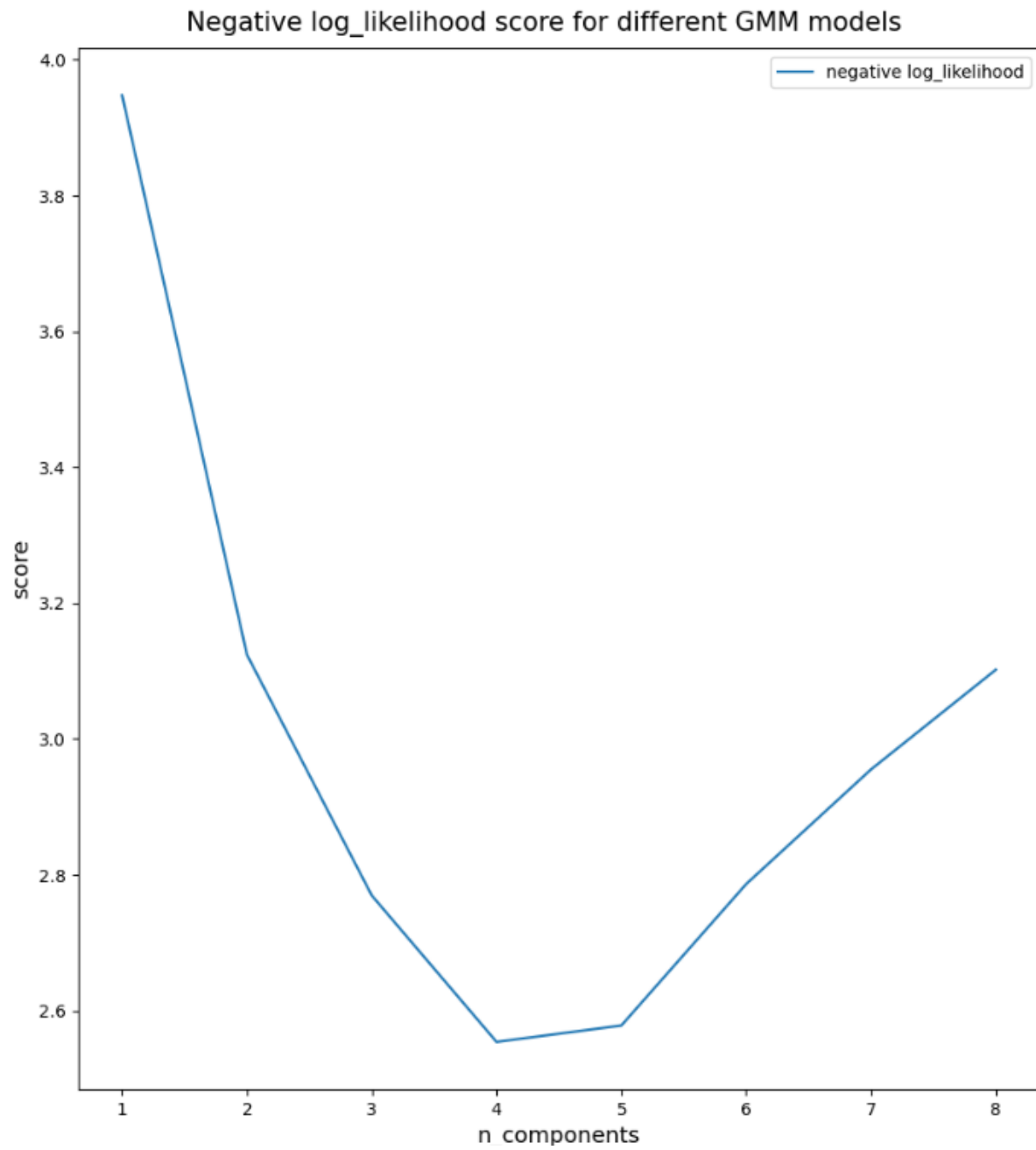
<b>GMM order selection based on (-) log-likelihood</b>	<b>Frequency</b>
<b>4</b>	<b>26</b>
<b>5</b>	<b>4</b>

The (-ve) log-likelihood vs n\_component graphs for one iteration (of the 30 experimental cycles) of every dataset is plotted and shown below:

**Dataset with 10 points:**

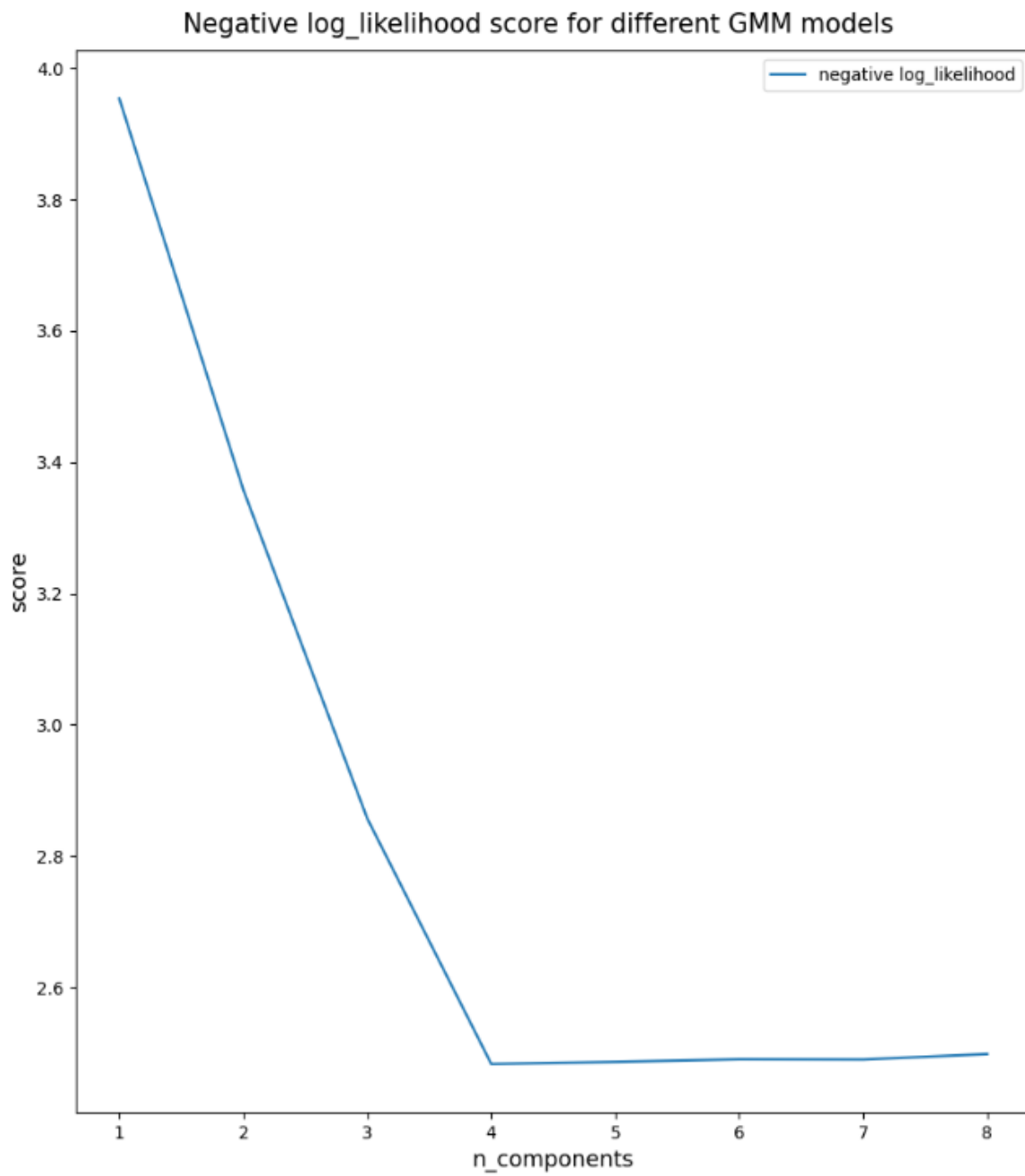


Dataset with 100 points:

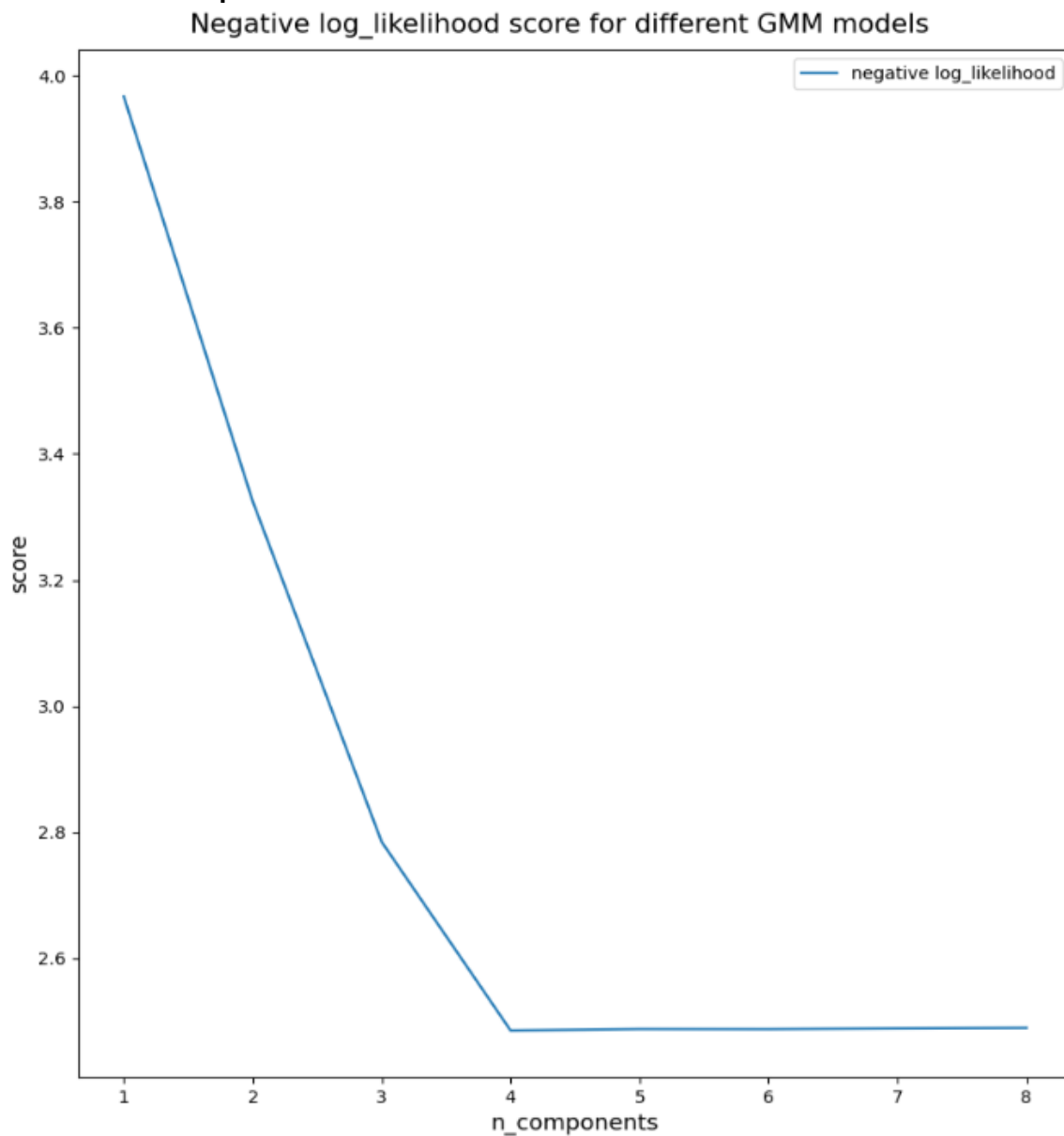




**Dataset with 1000 points:**



**Dataset with 10000 points:**



**Observations and Inference:**

- 1) For data sets with 100 points, 1000 points, and 10000 points the order that is mostly chosen based on minima of -ve likelihood score is 4. There are few occurrences of 5 and 6 but the most dominant order is 4.
- 2) The -ve likelihood vs n\_component graphs show which order produces the least -ve likelihood score and will be chosen accordingly. (best)

