# Group 28 - Report for Assignment 2 - Classification Models
## Pravin Ravi CE22B092
## Abhiram Sulige AE23B059
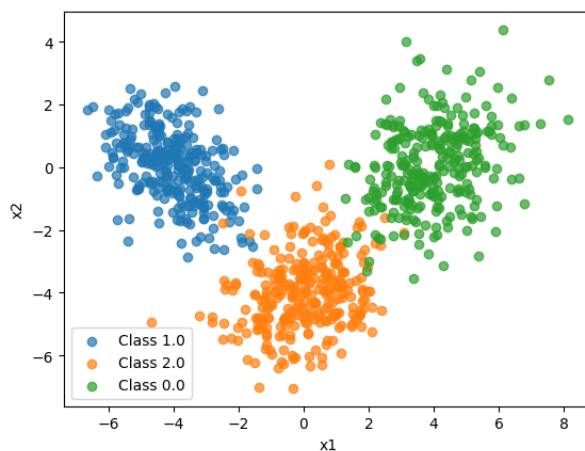
# Contents

# 1 Dataset 1

This is a 3-class 2 dimensional linearly separable dataset.



## 1.1 K-Nearest Neighbours Classifier

| K | Train Accuracy | Validation Accuracy | Test Accuracy |
|---|---|---|---|
| 1 | 100.00 | 97.92 | 95.83 |
| 5 | 98.57 | 98.33 | 97.50 |
| 9 | 97.74 | 98.33 | 95.83 |

We see that the best fit for the data is obtained for k = 5 model, it gives the best performance (most accuracy in test data).

### 1.1.1 Decision Region of Best Model



### 1.1.2 Performance Metrics

Accuracy for best performing model (k = 5):

Train Data Accuracy: 98.57%
Test Data Accuracy: 97.50%

Confusion Matrix:



Train Dataset



Test Dataset

For the test dataset predictions:

- Class 0

  - Precision: 1.000
  - Recall: 0.927
  - F1 Score: 0.962

- Class 1

  - Precision: 1.000
  - Recall: 1.000
  - F1 Score: 1.000

- Class 2

  - Precision: 0.925
  - Recall: 1.000
  - F1 Score: 0.961

**Average Precision: 0.975**
**Average Recall: 0.976**
**Average F1 Score: 0.974**

## 1.2 Bayes Classifier - Gaussian

### 1.2.1 Accuracy of each model

| Type | Train Accuracy | Validation Accuracy | Test Accuracy |
|------|----------------|---------------------|---------------|
| Same Cov. Matrix | 98.10 | 98.33 | 95.00 |
| Diff. Cov. Matrix | 98.21 | 98.33 | 95.00 |

The model with different covariance matrices for different classes has a better fit on training data, although performance on validation and test datasets are the same.

### 1.2.2 Decision Regions

**Same Covariance Matrix for all classes:**



The decision boundaries are linear.

**Different Covariance Matrix for each class:**



The decision boundaries are quadratic.

### 1.2.3 Performance Metrics

#### 1.2.3.1 a. Same Cov. Matrix for all classes   Confusion Matrix:



Train Dataset



Test Dataset

For the test dataset:

- **Class 0**
    - Precision: 1.000
    - Recall: 0.902
    - F1 Score: 0.949

- **Class 1**
    - Precision: 1.000
    - Recall: 0.952
    - F1 Score: 0.976

- **Class 2**
    - Precision: 0.860
    - Recall: 1.000
    - F1 Score: 0.925

**Average Precision: 0.953**
**Average Recall: 0.952**
**Average F1 Score: 0.950**

### 1.2.3.2 b. Diff Cov. Matrix for each class   Confusion Matrix:



Train Dataset



Test Dataset

For the test dataset:

- **Class 0**

    - Precision: 0.974
    - Recall: 0.902
    - F1 Score: 0.937

- **Class 1**

    - Precision: 1.000
    - Recall: 0.976
    - F1 Score: 0.988

- **Class 2**

    - Precision: 0.878
    - Recall: 0.973
    - F1 Score: 0.923

**Average Precision: 0.951**
**Average Recall: 0.951**
**Average F1 Score: 0.949**

# 2 Dataset 2

2-Class 2 Dimensional Non-linearly separable dataset



## 2.1 K-Nearest Neighbours Classifier

| K | Train Accuracy | Validation Accuracy | Test Accuracy |
|---|---|---|---|
| 1 | 100.00 | 94.34 | 93.75 |
| 5 | 94.12 | 95.60 | 100.00 |
| 9 | 94.30 | 96.23 | 97.50 |

We see that the best fit for the data is obtained for k = 9 model, it gives the best performance (most accuracy in test data).

### 2.1.1 Decision Region of Best Model



### 2.1.2 Performance Metrics

Accuracy for best performing model (k = 9):

Train Data Accuracy: 94.30%
Test Data Accuracy: 97.50%

Confusion Matrix:



Train Dataset



Test Dataset

For the test dataset predictions:

- **Class 0**

  - Precision: 0.974
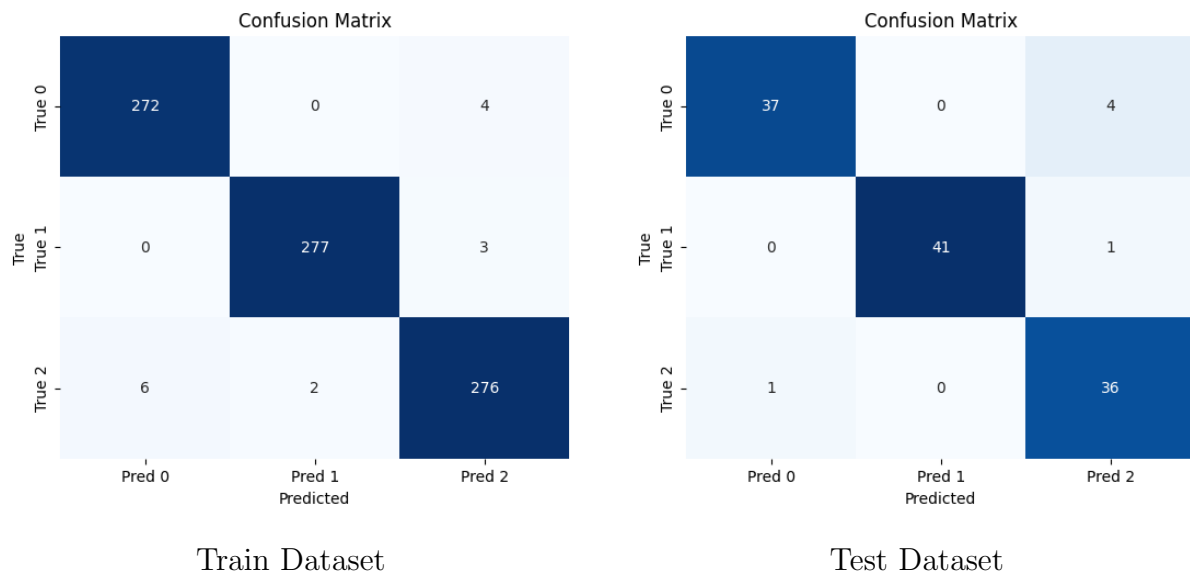  - Recall: 0.974
  - F1 Score: 0.974

- **Class 1**

  - Precision: 0.976
  - Recall: 0.976
  - F1 Score: 0.976

**Average Precision: 0.975**
**Average Recall: 0.975**
**Average F1 Score: 0.975**

## 2.2 Bayes Classifier - Gaussian

### 2.2.1 Accuracy of each model

| Type | Train Accuracy | Validation Accuracy | Test Accuracy |
|---|---|---|---|
| Same Cov. Matrix | 85.03 | 85.54 | 86.25 |
| Diff. Cov. Matrix | 84.85 | 84.91 | 86.25 |

### 2.2.2 Decision Regions

**Same Covariance Matrix for all classes:**



The decision boundaries are linear.

**Different Covariance Matrix for each class:**



The decision boundaries are quadratic, still looks very similar to above case.

### 2.2.3    Performance Metrics

#### 2.2.3.1    a. Same Cov. Matrix for all classes    Confusion Matrix:



Train Dataset



Test Dataset

For the test dataset:

- **Class 0**

  - Precision: 0.909
  - Recall: 0.789
  - F1 Score: 0.845

- **Class 1**
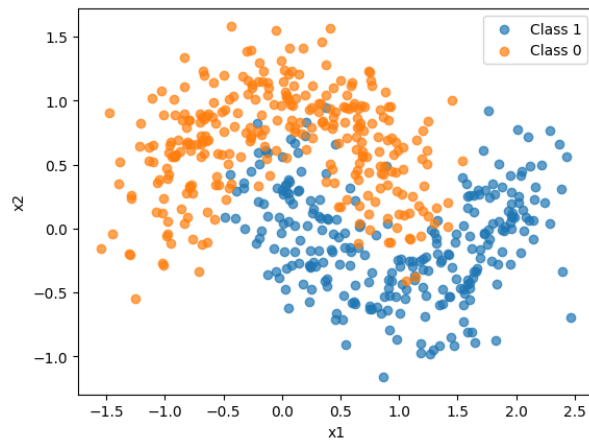
  - Precision: 0.830
  - Recall: 0.929
  - F1 Score: 0.876

**Average Precision: 0.869**
**Average Recall: 0.859**
**Average F1 Score: 0.861**

#### 2.2.3.2    b. Diff Cov. Matrix for each class    Confusion Matrix:

Train Dataset



Test Dataset

For the test dataset:

- **Class 0**

  - Precision: 0.909
  - Recall: 0.789
  - F1 Score: 0.845

- **Class 1**

  - Precision: 0.830
  - Recall: 0.929
  - F1 Score: 0.876

**Average Precision: 0.869**
**Average Recall: 0.859**
**Average F1 Score: 0.861**
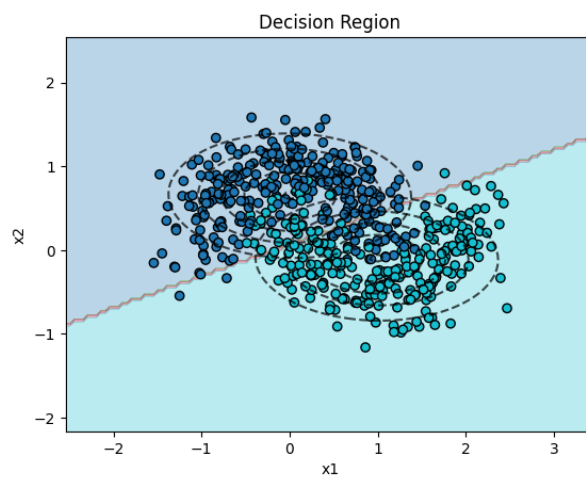
## 2.3 Naive-Bayes classifier with a Gaussian distribution for every class

### 2.3.1 Accuracy for same covariance matrices for all the classes

| Train Accuracy | Validation Accuracy | Test Accuracy |
|:---:|:---:|:---:|
| 84.5 | 85.5 | 86.2 |

### 2.3.2 Accuracy for different covariance matrices

| Train Accuracy | Validation Accuracy | Test Accuracy |
|:---:|:---:|:---:|
| 85.0 | 85.5 | 86.2 |

We see that both the models perform equally good, but the model with different covariance matrices for different classes has a better fit on training data. So let us choose the model with different covariance matrices for different classes as the best model.

### 2.3.3 Decision Region of Best Model



### 2.3.4 Performance Metrics

Accuracy for best performing model (different covariance matrices):

Train Data Accuracy: 85%
Validation Data Accuracy: 85.5%
Test Data Accuracy: 86.2%

Confusion Matrix:



Train Dataset



Test Dataset

For the test dataset predictions:

- Class 0

    - Precision: 0.91
    - Recall: 0.79
    - F1 Score: 0.84

- Class 1
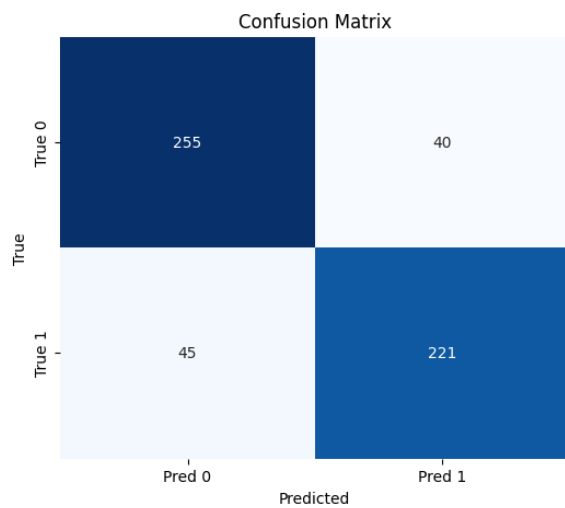
    - Precision: 0.83
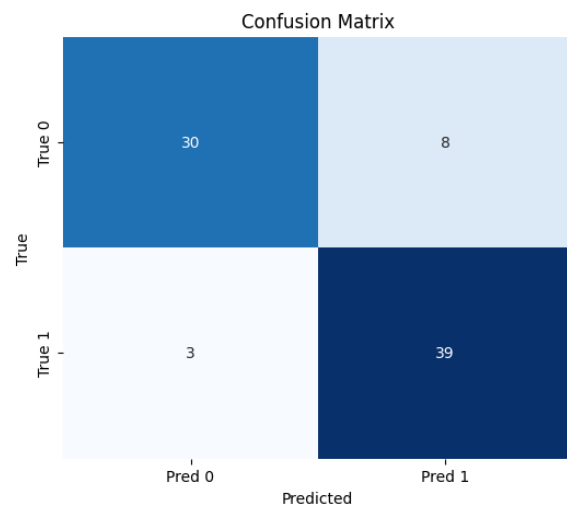    - Recall: 0.86
    - F1 Score: 0.86

**Average Precision: 0.87**
**Average Recall: 0.86**
**Average F1 Score: 0.86**

## 2.4  GMM based classifier

### 2.4.1  Accuracy with Q = 4

| Covariance_type | Train Accuracy | Validation Accuracy | Test Accuracy |
| --- | --- | --- | --- |
| Full cov | 93.4 | 95.0 | 97.5 |
| Diag cov | 92.5 | 94.3 | 97.5 |

### 2.4.2  Accuracy with Q = 6

| Covariance_type | Train Accuracy | Validation Accuracy | Test Accuracy |
| --- | --- | --- | --- |
| Full cov | 93.9 | 95.0 | 97.5 |
| Diag cov | 92.5 | 94.3 | 97.5 |

### 2.4.3  Accuracy with Q = 8

| Covariance_type | Train Accuracy | Validation Accuracy | Test Accuracy |
| --- | --- | --- | --- |
| Full cov | 93.8 | 94.3 | 97.5 |
| Diag cov | 92.5 | 94.3 | 97.5 |

### 2.4.4  Accuracy with Q = 10

| Covariance_type | Train Accuracy | Validation Accuracy | Test Accuracy |
| --- | --- | --- | --- |
| Full cov | 93.8 | 95.0 | 97.5 |
| Diag cov | 92.5 | 94.3 | 97.5 |

We see that the best fit for the data is obtained for Q = 6 (covariance_type = 'full') model, although all the model perform more or less equally better, this model is slightly best among them.

### 2.4.5 Decision Region of Best Model



### 2.4.6 Performance Metrics

Accuracy for best performing model:

Train Data Accuracy: 93.9%
Validation Data Accuracy: 95.0%
Test Data Accuracy: 97.5%

Confusion Matrix:



Train Dataset



Test Dataset

For the test dataset predictions:

- Class 0

- Precision: 0.974
- Recall: 0.974
- F1 Score: 0.974

- Class 1
  - Precision: 0.976
  - Recall: 0.976
  - F1 Score: 0.976

**Average Precision: 0.975**
**Average Recall: 0.975**
**Average F1 Score: 0.975**

# 3 Dataset 3

Image data set (Dimension of feature vector: 36) for 5 classes

## 3.1 K-Nearest Neighbours Classifier

| K | Train Accuracy | Validation Accuracy | Test Accuracy |
|---|---|---|---|
| 1 | 100.00 | 41.40 | 44.40 |
| 9 | 66.95 | 45.00 | 49.80 |
| 15 | 61.50 | 48.60 | 49.40 |

We see that the best fit for the data is obtained for k = 15 model, it gives the best performance (most accuracy in test data).
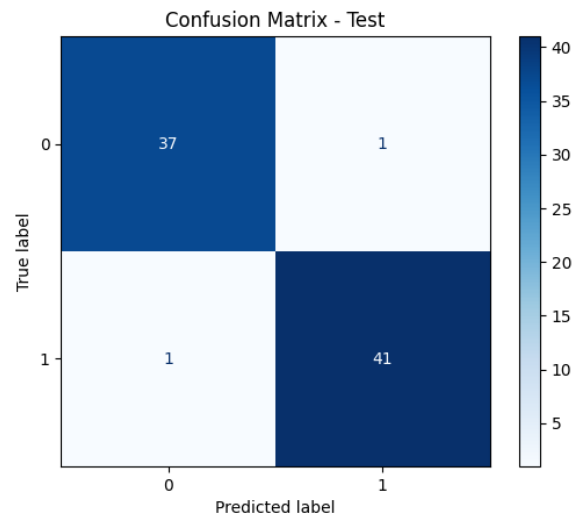
### 3.1.1 Performance Metrics

Accuracy for best performing model (k = 15):

$$\text{Train Data Accuracy: } 48.60\%$$
$$\text{Test Data Accuracy: } 49.40\%$$

Confusion Matrix:



Train Dataset          Test Dataset

For the test dataset predictions:

- **Class 0**

    - Precision: 0.571
    - Recall: 0.400
    - F1 Score: 0.471

- **Class 1**
  - Precision: 0.446
  - Recall: 0.450
  - F1 Score: 0.448

- **Class 2**
  - Precision: 0.421
  - Recall: 0.640
  - F1 Score: 0.508

- **Class 3**
  - Precision: 0.816
  - Recall: 0.620
  - F1 Score: 0.705

- **Class 4**
  - Precision: 0.356
  - Recall: 0.360
  - F1 Score: 0.358

**Average Precision: 0.522**
**Average Recall: 0.494**
**Average F1 Score: 0.498**

## 3.2 Bayes Classifier - Gaussian

### 3.2.1 Accuracy of each model

| Type | Train Accuracy | Validation Accuracy | Test Accuracy |
|---|---|---|---|
| Same Cov. Matrix | 53.20 | 49.00 | 50.40 |
| Diff. Cov. Matrix | 62.45 | 41.80 | 41.00 |

### 3.2.2 Performance Metrics

#### 3.2.2.1 a. Same Cov. Matrix for all classes   Confusion Matrix:



Train Dataset



Test Dataset

For the test dataset:

- **Class 0**

    - Precision: 0.597
    - Recall: 0.460
    - F1 Score: 0.520

- **Class 1**

    - Precision: 0.491
    - Recall: 0.530
    - F1 Score: 0.510

- **Class 2**

    - Precision: 0.425
    - Recall: 0.450
    - F1 Score: 0.437

- **Class 3**

  – Precision: 0.688

  – Recall: 0.750

  – F1 Score: 0.718

- **Class 4**

  – Precision: 0.330

  – Recall: 0.330

  – F1 Score: 0.330

**Average Precision: 0.506**
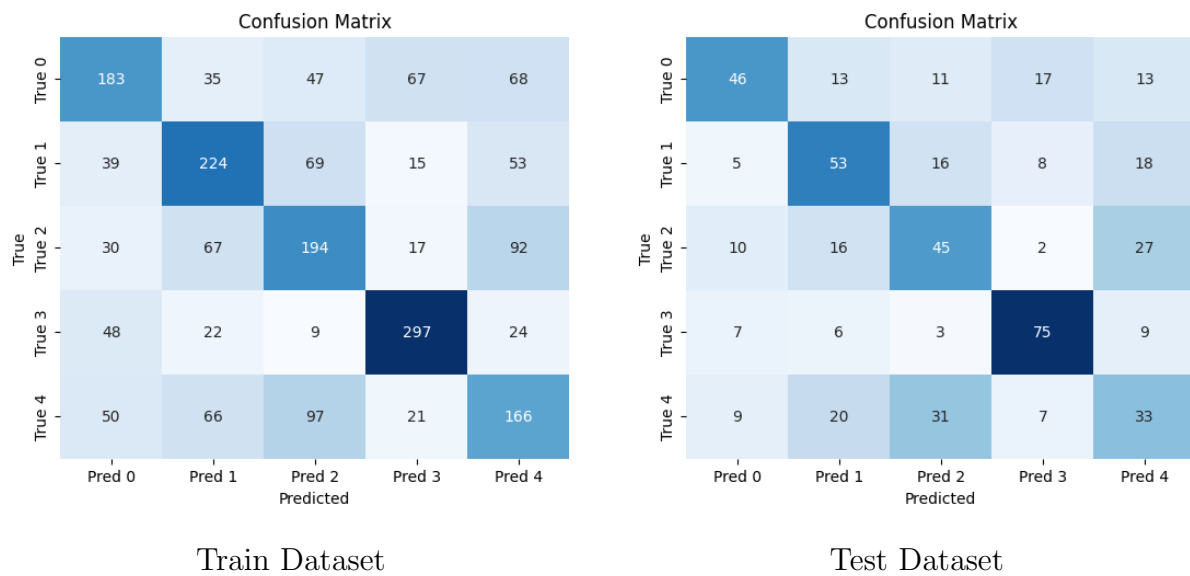**Average Recall: 0.504**
**Average F1 Score: 0.503**

### 3.2.2.2 b. Diff Cov. Matrix for each class  Confusion Matrix:



Train Dataset



Test Dataset

For the test dataset:

- **Class 0**

  – Precision: 0.321

  – Recall: 0.590

  – F1 Score: 0.415

- **Class 1**

  – Precision: 0.500

  – Recall: 0.220

  – F1 Score: 0.306

- **Class 2**

- Precision: 0.449
- Recall: 0.220
- F1 Score: 0.295

- **Class 3**

  - Precision: 0.481
  - Recall: 0.910
  - F1 Score: 0.630

- **Class 4**

  - Precision: 0.324
  - Recall: 0.110
  - F1 Score: 0.164

**Average Precision: 0.415**
**Average Recall: 0.410**
**Average F1 Score: 0.362**

## 3.3 Naive-Bayes classifier with a Gaussian distribution for every class

### 3.3.1 Accuracy for same covariance matrices for all the classes

| Train Accuracy | Validation Accuracy | Test Accuracy |
|:---:|:---:|:---:|
| 53.4 | 47.4 | 49.8 |

### 3.3.2 Accuracy for different covariance matrices

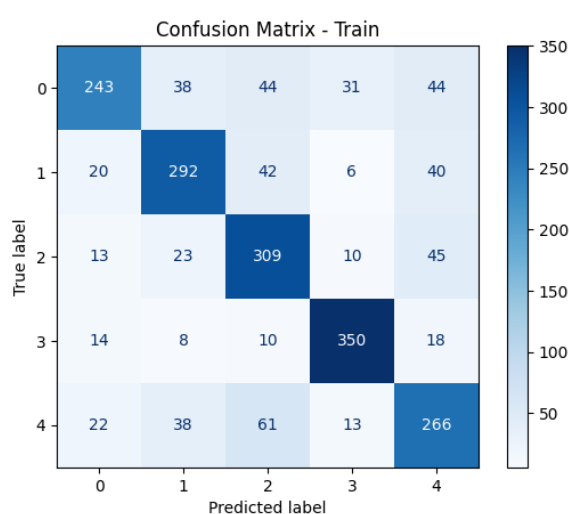| Train Accuracy | Validation Accuracy | Test Accuracy |
|:---:|:---:|:---:|
| 73.0 | 46.6 | 51.4 |

We see that model with different covariance matrices for different classes has a better accuracy on testing data with less difference in validation accuracy with that of other one, and also fits the training data better. So let us choose the model with different covariance matrices for different classes as the best model.
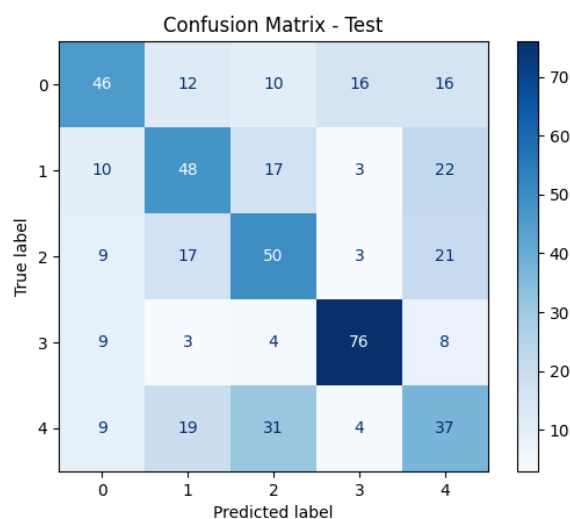
### 3.3.3 Performance Metrics

Accuracy for best performing model (different covariance matrices):

Train Data Accuracy: 73%
Validation Data Accuracy: 46.6%
Test Data Accuracy: 51.4%

Confusion Matrix:



Train Dataset

Test Dataset

For the test dataset predictions:

- Class 0

  - Precision: 0.554
  - Recall: 0.46
  - F1 Score: 0.50

- Class 1

  - Precision: 0.485
  - Recall: 0.48
  - F1 Score: 0.482

- Class 2

  - Precision: 0.446
  - Recall: 0.5
  - F1 Score: 0.472

- Class 3

  - Precision: 0.745
  - Recall: 0.76
  - F1 Score: 0.752

- Class 4

  - Precision: 0.356
  - Recall: 0.37
  - F1 Score: 0.363

**Average Precision: 0.52**
**Average Recall: 0.51**
**Average F1 Score: 0.51**

## 3.4 GMM based classifier

### 3.4.1 Accuracy with Q = 2

| Covariance_type | Train Accuracy | Validation Accuracy | Test Accuracy |
|---|---|---|---|
| Full cov | 85.2 | 47.2 | 49.0 |
| Diag cov | 54.4 | 49.6 | 46.0 |

### 3.4.2 Accuracy with Q = 3

| Covariance_type | Train Accuracy | Validation Accuracy | Test Accuracy |
|---|---|---|---|
| Full cov | 92.1 | 44.4 | 44.8 |
| Diag cov | 54.4 | 49.6 | 46.0 |

### 3.4.3 Accuracy with Q = 4

| Covariance_type | Train Accuracy | Validation Accuracy | Test Accuracy |
|---|---|---|---|
| Full cov | 95.4 | 45.8 | 44.6 |
| Diag cov | 54.4 | 49.6 | 46.0 |

### 3.4.4 Accuracy with Q = 5

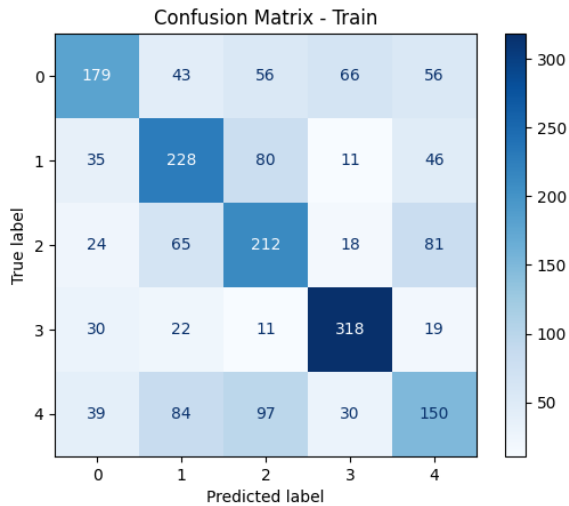| Covariance_type | Train Accuracy | Validation Accuracy | Test Accuracy |
|---|---|---|---|
| Full cov | 97.8 | 41.6 | 40.8 |
| Diag cov | 54.4 | 49.6 | 46.0 |

We see that all the models with (covariance type : diagonal) perform equally good (from validation data accuracy). Among them, we can choose any one, We choose the model for Q = 2 (covariance_type = 'diag').

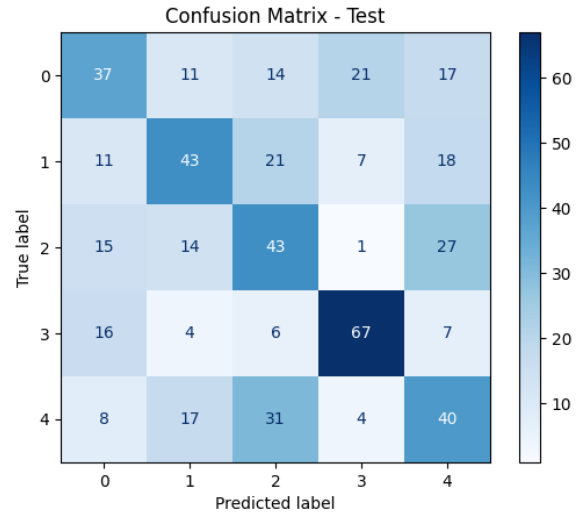### 3.4.5 Performance Metrics

Accuracy for best performing model:

Train Data Accuracy: 54.4%
Validation Data Accuracy: 49.6%
Test Data Accuracy: 46%

Confusion Matrix:

Train Dataset



Test Dataset

For the test dataset predictions:

- Class 0

  - Precision: 0.425
  - Recall: 0.37
  - F1 Score: 0.96

- Class 1

  - Precision: 0.483
  - Recall: 0.43
  - F1 Score: 0.455

- Class 2

  - Precision: 0.374
  - Recall: 0.43
  - F1 Score: 0.4

- Class 3

  - Precision: 0.67
  - Recall: 0.67
  - F1 Score: 0.67

- Class 4

  - Precision: 0.37
  - Recall: 0.4
  - F1 Score: 0.38

**Average Precision: 0.464**
**Average Recall: 0.46**
**Average F1 Score: 0.461**