

# Introduction to Machine Learning

## HOMEWORK 1

AKSHATA KUMBLE

Q1.

PART A:

The pdf for  $X$  is given as a mixture:

$$p(x) = p(x|L=0)P(L=0) + p(x|L=1)P(L=1)$$

where  $L$  is the true class label (0 or 1) and the class priors are

$$P(L=0) = 0.35 \text{ and } P(L=1) = 0.65$$

The minimum expected risk classification rule is based on a likelihood ratio test for two classes using gaussian distributions

$$\frac{p(x|L=1)}{p(x|L=0)} \geq r$$

where the threshold  $r$  is determined by:

$$r = \frac{P(L=0)}{P(L=1)} \cdot \frac{\lambda_{10}}{\lambda_{01}}$$

$p(x|L=1)$  is the gaussian pdf for class 1 (with mean  $m_1$  & covariance  $C_1$ )

$p(x|L=0)$  is the gaussian pdf for class 0 (with mean  $m_0$  & covariance  $C_0$ )

$r$  depends on class priors & loss values

Gaussian pdf is given by:

$$p(x|m, C) = \frac{1}{\sqrt{(2\pi)^d |C|}} e^{-\frac{1}{2}(x-m)^T C^{-1} (x-m)}$$

Given:

Class priors:  $P(L=0) = 0.35$ ,  $P(L=1) = 0.65$

Loss values: Assume symmetric loss  $\rightarrow \lambda_{10} = \lambda_{01} = 1$

$\therefore$  Threshold  $r$  becomes:

$$r = \frac{0.35}{0.65} \approx 0.538$$

# Implement the classifier

I generated 10,000 samples from the two gaussian distributions for each class:

Class 0 has mean vector  $m_0 = [-1, -1, -1, -1]$  and covariance matrix  $C_0$ .

Class 1 has mean vector  $m_1 = [1, 1, 1, 1]$  and covariance matrix  $C_1$ .

Sample generation using multivariate gaussian distribution to generate samples for each class.

```
In [1]: 1 import numpy as np
2
3 # Parameters
4 m0 = np.array([-1, -1, -1, -1])
5 C0 = np.array([[2, -0.5, 0.3, 0],
6                 [-0.5, 1, -0.5, 0],
7                 [0.3, -0.5, 1, 0],
8                 [0, 0, 0, 2]])
9
10 m1 = np.array([1, 1, 1, 1])
11 C1 = np.array([[1, 0.3, -0.2, 0],
12                  [0.3, 2, 0.3, 0],
13                  [-0.2, 0.3, 1, 0],
14                  [0, 0, 0, 3]])
15
16 N = 10000 # Number of samples
17
18 # Class priors
19 P0 = 0.35
20 P1 = 0.65
21
22 # Generate samples
23 samples_class0 = np.random.multivariate_normal(m0, C0, int(N * P0))
24 samples_class1 = np.random.multivariate_normal(m1, C1, int(N * P1))
25
26 # Combine samples and labels
27 samples = np.vstack((samples_class0, samples_class1))
28 labels = np.hstack((np.zeros(len(samples_class0)), np.ones(len(samples_class1))))
```

Computing likelihood ratio for each sample

```
In [2]: 1 from scipy.stats import multivariate_normal
2
3 def likelihood_ratio(x, m0, C0, m1, C1):
4     p_x_given_L0 = multivariate_normal.pdf(x, mean=m0, cov=C0)
5     p_x_given_L1 = multivariate_normal.pdf(x, mean=m1, cov=C1)
6     return p_x_given_L1 / p_x_given_L0
7
```

Varying the threshold  $\gamma$  from 0 to  $\infty$

True positive rate:  $P(D=1 | L=1)$

False positive rate:  $P(D=1 | L=0)$

```
In [3]: 1 gammas = np.logspace(-3, 3, 100) # Vary threshold
2
3 tpr = [] # True Positive Rate
4 fpr = [] # False Positive Rate
5
6 for gamma in gammas:
7     decisions = np.array([1 if likelihood_ratio(x, m0, C0, m1, C1) > gamma else 0 for x in samples])
8     tp = np.sum((decisions == 1) & (labels == 1)) / np.sum(labels == 1)
9     fp = np.sum((decisions == 1) & (labels == 0)) / np.sum(labels == 0)
10    tpr.append(tp)
11    fpr.append(fp)
12
```

# ROC and minimum probability of error

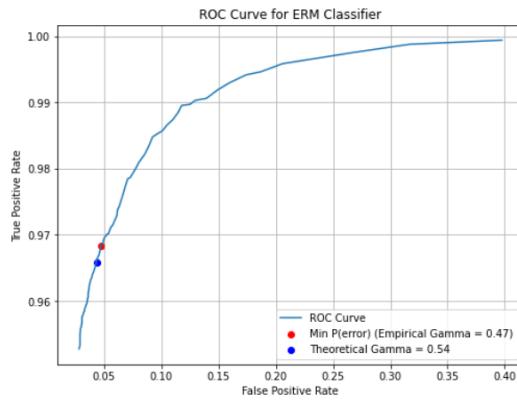
```
In [4]: M
1 import matplotlib.pyplot as plt
2
3 plt.plot(fpr, tpr, label="ROC Curve")
4 plt.xlabel("False Positive Rate")
5 plt.ylabel("True Positive Rate")
6 plt.title("ROC Curve for ERM Classifier")
7 plt.show()
8
9
10 errors = [fp * P0 + (1 - tp) * P1 for tp, fp in zip(tpr, fpr)]
11 min_error_gamma = gammas[np.argmin(errors)]
12 min_error = min(errors)
13
14 print(f"Minimum Probability of Error: {min_error} at gamma = {min_error_gamma}")
```

The min probability of error occurs when

$$P(\text{error}; \gamma) = P(D=1 | L=0; \gamma)P(L=0) + P(D=0 | L=1; \gamma)P(L=1)$$

is minimized

Theoretical gamma: 0.5384615384615384



Empirical Minimum Probability of Error: 0.03709999999999994 at gamma = 0.4700000000000003  
Error at Theoretical Gamma = 0.5384615384615384: 0.03740000000000003

(0.47)

The empirical  $\gamma$  is slightly lower than (0.53)  
the theoretical  $\gamma$  likely due to the finite sample size and stochastic nature of the generated data. The difference in the empirical min. error & error at theoretical gamma is very small, which suggests that the classifier performs well near the theoretical optimal point.

## PART B.

### Naive Bayesian Classifier

Assume that features are independent given each class label. That means all the off-diagonal elements in the covariance matrix are set to 0

#### Diagonal Covariance Matrices:

- For class 0, the covariance matrix  $C_0$  is modified to have only the diagonal elements from the original  $C_0$ :

$$C_0^{\text{NB}} = \text{diag}(2, 1, 1, 2)$$

- For class 1, covariance matrix  $C_1$  is:

$$C_1^{\text{NB}} = \text{diag}(1, 2, 1, 3)$$

#### Likelihood Computation:

In Naive Bayes setting, the gaussian pdf simplifies since we only consider the diagonal elements of the covariance matrix.

The pdfs for both classes  $L=0, L=1$  are computed as:

$$p(x|L=j) = \prod_{i=1}^4 \frac{1}{\sqrt{2\pi\sigma_{ji}^2}} \exp\left(-\frac{(x_i - m_{ji})^2}{2\sigma_{ji}^2}\right)$$

$\sigma_{ji}^2$  is variance for  $i^{\text{th}}$  feature in class  $L=j$

$m_{ji}$  is mean for  $i^{\text{th}}$  feature in class  $L=j$

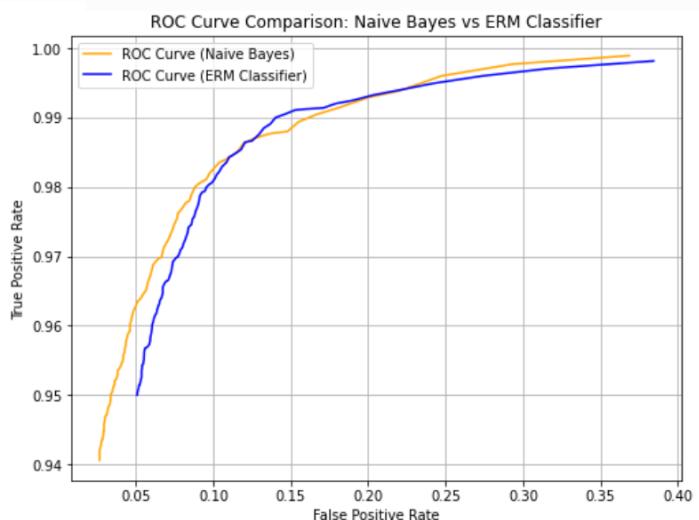
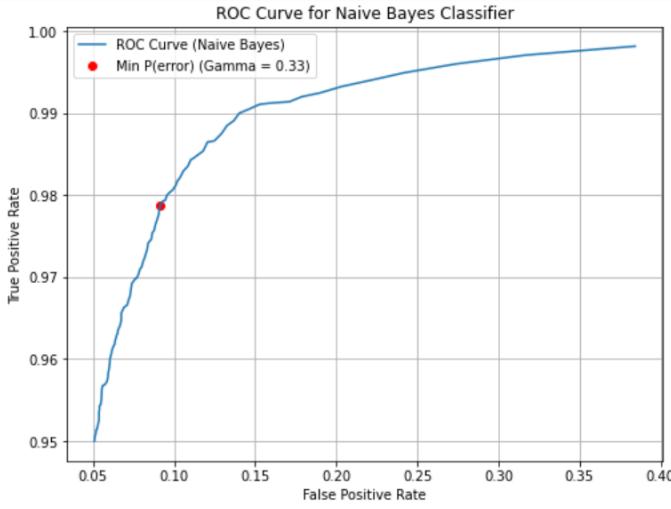
Naive Bayes produces lower performing ROC curve. It is a less accurate model because it ignores feature correlation.

This model mismatch increases the error rate.

```

In [16]: M
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.stats import multivariate_normal
4
5 # Parameters (Same means)
6 m0 = np.array([-1, -1, -1, -1])
7 m1 = np.array([1, 1, 1, 1])
8
9 # Naive Bayes Assumption: Diagonal Covariance Matrices
10 C0_naive = np.diag(np.diag([[2, -0.5, 0.3, 0],
11                             [-0.5, 1, -0.5, 0],
12                             [0.3, -0.5, 1, 0],
13                             [0, 0, 0, 2]]))
14
15 C1_naive = np.diag(np.diag([[1, 0.3, -0.2, 0],
16                             [0.3, 2, 0.3, 0],
17                             [-0.2, 0.3, 1, 0],
18                             [0, 0, 0, 3]])))
19
20 # Class priors
21 P0 = 0.35
22 P1 = 0.65
23
24 # Generate samples (Same as before)
25 samples_class0 = np.random.multivariate_normal(m0, C0_naive, int(N * P0))
26 samples_class1 = np.random.multivariate_normal(m1, C1_naive, int(N * P1))
27
28 # Combine samples and labels
29 samples = np.vstack((samples_class0, samples_class1))
30 labels = np.hstack((np.zeros(len(samples_class0)), np.ones(len(samples_class1))))
31
32 def likelihood_ratio_naive(x, m0, C0, m1, C1):
33     p_x_given_L0 = multivariate_normal.pdf(x, mean=m0, cov=C0)
34     p_x_given_L1 = multivariate_normal.pdf(x, mean=m1, cov=C1)
35     return p_x_given_L1 / p_x_given_L0
36
37 # Generate finer gamma range around expected values
38 gammas = np.linspace(0.01, 1, 100) # Vary gamma finely around expected range
39
40 tpr = [] # True Positive Rate
41 fpr = [] # False Positive Rate
42
43 for gamma in gammas:
44     decisions = np.array([1 if likelihood_ratio_naive(x, m0, C0_naive, m1, C1_naive) > gamma else 0 for x in samples])
45     tp = np.sum((decisions == 1) & (labels == 1)) / np.sum(labels == 1)
46     fp = np.sum((decisions == 1) & (labels == 0)) / np.sum(labels == 0)
47     tpr.append(tp)
48     fpr.append(fp)
49
50 # Plot ROC curve
51 plt.figure(figsize=(8, 6))
52 plt.plot(fpr, tpr, label="ROC Curve (Naive Bayes)")
53 plt.xlabel("False Positive Rate")
54 plt.ylabel("True Positive Rate")
55 plt.title("ROC Curve for Naive Bayes Classifier")
56 plt.grid(True)
57
58 # Compute error for each threshold gamma
59 errors = [fp * P0 + (1 - tp) * P1 for tp, fp in zip(tpr, fpr)]
60 min_error_idx = np.argmin(errors)
61 min_error_gamma = gammas[min_error_idx]
62 min_error = min(errors)
63
64 # Mark the minimum probability of error on the ROC curve
65 plt.scatter(fpr[min_error_idx], tpr[min_error_idx], color='red', label=f'Min P(error) (Gamma = {min_error_gamma:.2f})')
66 plt.legend()
67 plt.show()
68
69 # Report the minimum probability of error for Naive Bayes
70 print(f"Naive Bayes Minimum Probability of Error: {min_error} at gamma = {min_error_gamma}")
71

```



Naive Bayes Minimum Probability of Error: 0.0458000000000003 at gamma = 0.33

The ROC curve for naive bayes classifier is slightly worse than the true model's ROC curve. It is shifted slightly down and to the right compared to the true model. The model mismatch does negatively impact the ROC curve & the min. prob. error.

The Naive Bayes ROC is slightly inferior to the true models. At any given false +ve rate, the true +ve rate for Naive Bayes is slightly lower than for the true model.

The impact is noticeable but not severe suggesting that while there are correlations between features in the true distribution, they are not extremely strong.

The min. prob. error increased from 3.71% (true model) to 4.58% (Naive Bayes). This is an increase in error of about 23.45% relative to the true model's performance.

The shift in optimal gamma suggests that using the theoretical threshold based on priors might not be optimal for Naive Bayes classifier.

## PART C

### Fisher Linear Discriminant Analysis

LDA seeks a projection of the data onto a lower dimension space that maximizes the separation between the 2 classes

Using 10,000 samples the mean vectors and covariance matrices are estimated for each class

Estimated mean for class 0:

$$\hat{m}_0 = \frac{1}{N_0} \sum_{i=1}^{N_0} x_i$$

Estimated covariance for class 0:

$$\hat{C}_0 = \frac{1}{N_0 - 1} \sum_{i=1}^{N_0} (x_i - \hat{m}_0)(x_i - \hat{m}_0)^T$$

Similarly for Class 1.

Fisher LDA projection vector  $w_{LDA}$  is:

$$w_{LDA} = S_w^{-1} (m_1 - m_0)$$

where  $S_w = C_0 + C_1$

Then I project the samples onto the LDA vector:  $z = w_{LDA}^T x$

Classify the samples based on whether the projection  $z$  is greater than or less than some threshold  $\tau$

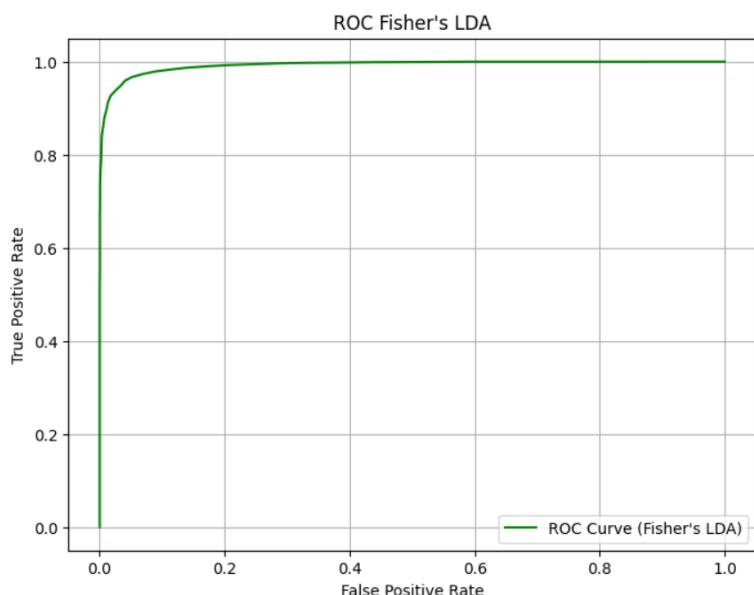
Then the threshold  $\tau$  is varied from  $-\infty$  to  $\infty$  & the ROC curve is plotted

Fisher LDA classifier often performs well when the data is linearly separable & produces better results than Naive Bayes.

```

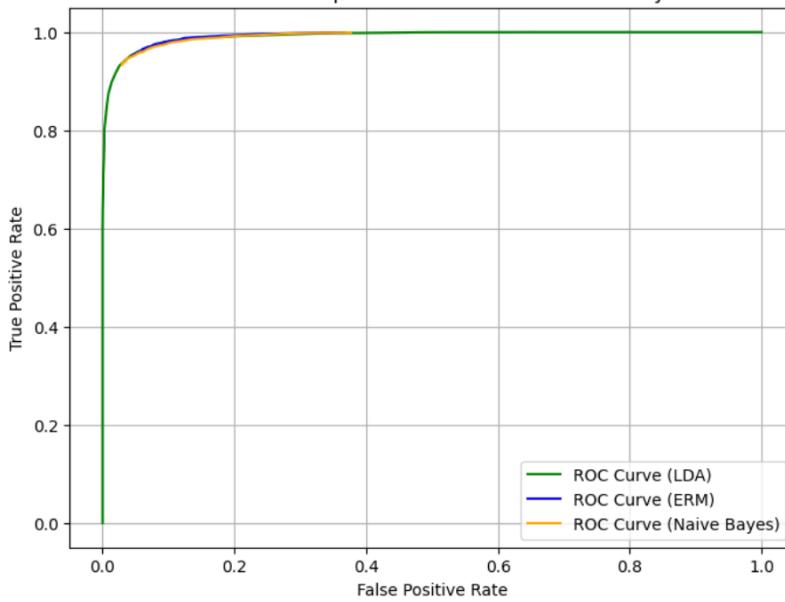
In [1]: 
 1 import numpy as np
 2 import matplotlib.pyplot as plt
 3
 4 # Parameters for the class means and covariance matrices
 5 m0 = np.array([-1, -1, -1, -1])
 6 m1 = np.array([1, 1, 1, 1])
 7
 8 C0 = np.array([[2, -0.5, 0.3, 0],
 9                 [-0.5, 1, -0.5, 0],
10                 [0.3, -0.5, 1, 0],
11                 [0, 0, 0, 2]])
12
13 C1 = np.array([[1, 0.3, -0.2, 0],
14                 [0.3, 2, 0.3, 0],
15                 [-0.2, 0.3, 1, 0],
16                 [0, 0, 0, 3]])
17
18 # Diagonal covariance matrices for Naive Bayes assumption
19 C0_diag = np.diag(C0)
20 C1_diag = np.diag(C1)
21
22 # Class priors
23 P0 = 0.35
24 P1 = 0.65
25
26 # Generate samples (from the original Gaussian distributions)
27 samples_class0 = np.random.multivariate_normal(m0, C0, int(10000 * P0))
28 samples_class1 = np.random.multivariate_normal(m1, C1, int(10000 * P1))
29
30 # Combine samples and labels
31 samples = np.vstack((samples_class0, samples_class1))
32 labels = np.hstack((np.zeros(len(samples_class0)), np.ones(len(samples_class1))))
33
34 # Define Fisher's LDA projection function
35 def fisher_lda(m0, C0, m1, C1):
36     Sw = C0 + C1 # Within-class scatter matrix
37     w_lda = np.linalg.inv(Sw).dot(m1 - m0) # Fisher LDA vector
38     return w_lda
39
40 # Get Fisher's LDA projection vector
41 w_lda = fisher_lda(m0, C0, m1, C1)
42
43 # Project the samples onto the LDA axis
44 projected_data = samples.dot(w_lda)
45
46 # Compute ROC curve for Fisher's LDA
47 gammas_lda = np.linspace(np.min(projected_data), np.max(projected_data), 100)
48 tpr_lda = [] # True Positive Rate for LDA
49 fpr_lda = [] # False Positive Rate for LDA
50
51 for gamma in gammas_lda:
52     decisions_lda = np.where(projected_data > gamma, 1, 0)
53     tp = np.sum((decisions_lda == 1) & (labels == 1)) / np.sum(labels == 1)
54     fp = np.sum((decisions_lda == 1) & (labels == 0)) / np.sum(labels == 0)
55     tpr_lda.append(tp)
56     fpr_lda.append(fp)
57
58 plt.figure(figsize=(8, 6))
59 plt.plot(fpr_lda, tpr_lda, label="ROC Curve (Fisher's LDA)", color='green')
60
61 plt.xlabel("False Positive Rate")
62 plt.ylabel("True Positive Rate")
63 plt.title("ROC Fisher's LDA")
64 plt.legend()
65 plt.grid(True)
66 plt.show()
67
68 print(f"Naive Bayes: {np.interp(fpr_naive, tpr_naive)}")
69 print(f"Fisher LDA: {np.interp(fpr_lda, tpr_lda)}")
70

```



Fisher's LDA Minimum Probability of Error: 0.0396000000000004 at gamma = 0.43

ROC Curve Comparison: LDA vs ERM vs Naive Bayes



Fisher's LDA Minimum Probability of Error: 0.04429999999999997 at gamma = 0.45  
 Naive Bayes Minimum Probability of Error: 0.0462999999999999 at gamma = 0.34  
 ERM Classifier Minimum Probability of Error: 0.04360000000000014 at gamma = 0.54

All 3 ROC curves are very close to each other, indicating similar performance. However on closer inspection we observe that ERM marginally outperforms LDA & Naive Bayes. The difference is subtle but reflects better class separation.

LDA is known for finding the linear combination of features that best separates the classes.

Naive Bayes assumes that all features are conditionally independent thereby resulting in a slightly higher error rate.

ERM directly minimizes the empirical error rate, fits the data more precisely as it has the lowest minimum prob. error.

The slightly different thresholds indicate how each model adjusts its decision boundary.

ERM's optimal gamma is higher suggesting a more conservative threshold.

Q2.

Use a mixture of 4 gaussian distributions to generate 10,000 samples  
Class priors are set as  $P(L=1) = 0.3$ ,  $P(L=2) = 0.3$ ,  $P(L=3) = 0.4$  for classes 1, 2, 3 respectively. Class 1 and Class 2 have their own gaussian distribution with mean  $\mu_1, \mu_2$  and covariance  $\Sigma_1, \Sigma_2$ . Class 3 is a mixture of two gaussians with equal weight & means  $\mu_{3a}, \mu_{3b}$  & covariances  $\Sigma_{3a}, \Sigma_{3b}$ . We need to define the Gaussian parameters such that the means of the gaussians are separated by 2 to 3 times the average standard deviation of the components in order to ensure a moderate amount of overlap between the class conditional distribution which is important for testing the classifier.

e.g. of parameters for the three classes :

class 1 :  $\mu_1 = [0, 0, 0]$ ,  $\Sigma_1 = \text{diag}(1, 1, 1)$ ; class 2 :  $\mu_2 = [3, 3, 3]$ ,  $\Sigma_2 = \text{diag}(1, 1, 1)$

class 3 : Mixture of  $\mu_{3a} = [0, 5, 0]$ ,  $\Sigma_{3a} = \text{diag}(1, 1, 1)$ ;  $\mu_{3b} = [5, 0, 5]$ ,  $\Sigma_{3b} = \text{diag}(1, 1, 1)$

We need to :

1. Generate 10,000 samples
2. Specify a decision rule that achieves minimum probability error using 0-1 loss function

The classifier minimizes the probability of error, which classifies a sample a sample by maximizing the posterior probability .

$$\hat{L} = \arg \max L = \arg \max p(x|L) P(L)$$

where :  $p(x|L)$  is the class conditional pdf and  $P(L)$  are class priors

Computing the posterior probabilities for each class as :

$$P(L=1|x) = p(x|L=1) P(L=1), P(L=2|x) = p(x|L=2) P(L=2), P(L=3|x) = p(x|L=3) P(L=3)$$

Then we classify the sample to the class with the highest posterior probability given the sample

The decision rule for MAP classifier chooses the class with the highest posterior probability.

3. Implement the classifier and classify the samples by Computing the confusion matrix (which is a matrix that shows how many samples were correctly or incorrectly classified for each class).

In [1]:

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from mpl_toolkits.mplot3d import Axes3D
4 from scipy.stats import multivariate_normal
5
6 # Set random seed for reproducibility
7 np.random.seed(42)
8
9 def generate_data(n_samples=10000):
10     # Define parameters for the Gaussian distributions
11     means = [
12         np.array([0, 0, 0]),      # Class 1
13         np.array([2.5, 2.5, 2.5]), # Class 2 - first component
14         np.array([-2.5, -2.5, -2.5]), # Class 2 - second component
15         np.array([2.5, -2.5, 2.5]) # Class 3
16     ]
17
18     cov = np.eye(3) # Standard deviation of 1 in each dimension
19
20     # Generating samples
21     X = np.zeros((n_samples, 3))
22     y = np.zeros(n_samples)
23
24     class_priors = [0.3, 0.3, 0.4] # Given in the question
25
26     current_idx = 0
27     for class_idx, prior in enumerate(class_priors):
28         n_class_samples = int(n_samples * prior)
29
30         if class_idx == 1: # Class 2 has two components
31             n_per_component = n_class_samples // 2
32             X[current_idx:current_idx + n_per_component] = np.random.multivariate_normal(
33                 means[1], cov, n_per_component)
34             X[current_idx + n_per_component:current_idx + n_class_samples] = np.random.multivariate_normal(
35                 means[2], cov, n_class_samples - n_per_component)
36         else:
37             X[current_idx:current_idx + n_class_samples] = np.random.multivariate_normal(
38                 means[class_idx if class_idx < 1 else 3], cov, n_class_samples)
39
40         y[current_idx:current_idx + n_class_samples] = class_idx + 1
41         current_idx += n_class_samples
42
43     return X, y
44
45 def calculate_class_probabilities(X, means, cov, class_priors):
46     probabilities = np.zeros((X.shape[0], 3))
47
48     # Class 1 and 3 have single Gaussian
49     probabilities[:, 0] = multivariate_normal.pdf(X, means[0], cov) * class_priors[0]
50     probabilities[:, 2] = multivariate_normal.pdf(X, means[3], cov) * class_priors[2]
51
52     # Class 2 has mixture of two Gaussians
53     class2_prob = (multivariate_normal.pdf(X, means[1], cov) +
54                     multivariate_normal.pdf(X, means[2], cov)) / 2
55     probabilities[:, 1] = class2_prob * class_priors[1]
56
57     return probabilities
58
59 def bayes_classifier(X, means, cov, class_priors):
60     probabilities = calculate_class_probabilities(X, means, cov, class_priors)
61     return np.argmax(probabilities, axis=1) + 1
62
63 def erm_classifier(X, means, cov, class_priors, loss_matrix):
64     probabilities = calculate_class_probabilities(X, means, cov, class_priors)
65
66     # Calculate expected risk for each decision
67     risks = np.zeros((X.shape[0], 3))
68     for i in range(3): # For each possible decision
69         for j in range(3): # For each true class
70             risks[:, i] += loss_matrix[i, j] * probabilities[:, j]
71
72     return np.argmin(risks, axis=1) + 1
73
74 def evaluate_classifier(y_true, y_pred):
75     confusion_matrix = np.zeros((3, 3))
76     for i in range(len(y_true)):
77         confusion_matrix[int(y_pred[i]-1), int(y_true[i]-1)] += 1
78
79     accuracy = np.sum(y_true == y_pred) / len(y_true)
80     return confusion_matrix, accuracy
81
82 def plot_results(X, y_true, y_pred):
83     fig = plt.figure(figsize=(10, 10))
84     ax = fig.add_subplot(111, projection='3d')
85
86     colors = ['g' if pred == true else 'r' for pred, true in zip(y_pred, y_true)]
87     markers = {1: 'o', 2: '^', 3: 's'}
88
89     for class_label in [1, 2, 3]:
90         mask = y_true == class_label
91         ax.scatter(X[mask, 0], X[mask, 1], X[mask, 2],
92                    c=[colors[i] for i in range(len(colors)) if y_true[i] == class_label],
93                    marker=markers[class_label], label=f'Class {class_label}')
94
95     ax.set_xlabel('X')
96     ax.set_ylabel('Y')
97     ax.set_zlabel('Z')
98     ax.legend()
99     plt.show()
100 X, y = generate_data()
101 # Defining parameters
102 means = [
103     np.array([0, 0, 0]),
104     np.array([2.5, 2.5, 2.5]),
105     np.array([-2.5, -2.5, -2.5]),
106     np.array([2.5, -2.5, 2.5])
107 ]
108 cov = np.eye(3)
109 class_priors = [0.3, 0.3, 0.4]
110

```

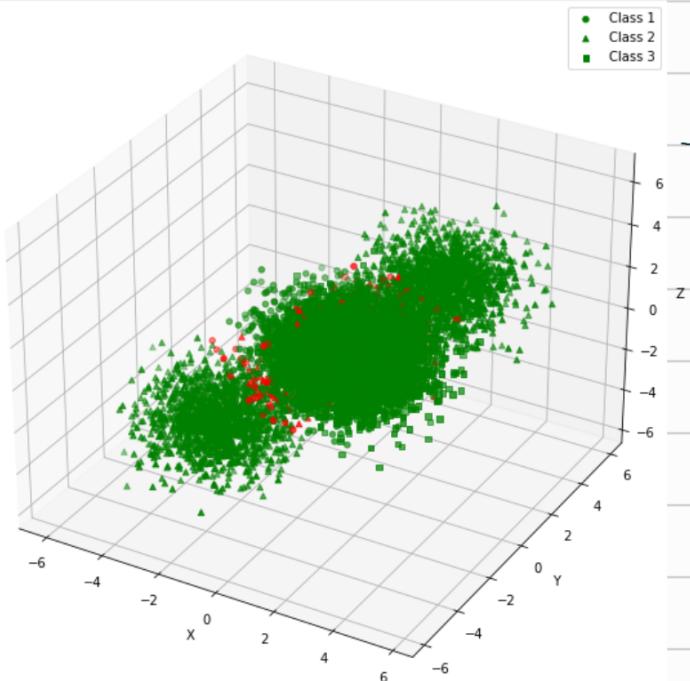
```

111 # Part A: Bayes Classifier
112 y_pred_bayes = bayes_classifier(X, means, cov, class_priors)
113 conf_matrix_bayes, accuracy_bayes = evaluate_classifier(y, y_pred_bayes)
114
115 print("Part A Results:")
116 print(f"Bayes Classifier Accuracy: {accuracy_bayes:.4f}")
117 print("Confusion Matrix:")
118 print(conf_matrix_bayes)
119
120 # Visualization
121 plot_results(X, y, y_pred_bayes)
122
123 # Part B: ERM Classifier with different loss matrices
124 loss_matrix_10 = np.array([[0, 10, 10],
125                           [1, 0, 10],
126                           [1, 1, 0]])
127
128 loss_matrix_100 = np.array([[0, 100, 100],
129                           [1, 0, 100],
130                           [1, 1, 0]])
131
132 y_pred_erm_10 = erm_classifier(X, means, cov, class_priors, loss_matrix_10)
133 y_pred_erm_100 = erm_classifier(X, means, cov, class_priors, loss_matrix_100)
134
135 _, accuracy_erm_10 = evaluate_classifier(y, y_pred_erm_10)
136 _, accuracy_erm_100 = evaluate_classifier(y, y_pred_erm_100)
137
138 print("\nPart B Results:")
139 print(f"ERM Classifier (Loss Matrix 10) Accuracy: {accuracy_erm_10:.4f}")
140 print(f"ERM Classifier (Loss Matrix 100) Accuracy: {accuracy_erm_100:.4f}")
141
142 # Plot results for ERM classifiers
143 plot_results(X, y, y_pred_erm_10)
144 plot_results(X, y, y_pred_erm_100)

```

## VISualization

Part A Results:  
 Bayes Classifier Accuracy: 0.9748  
 Confusion Matrix:  
 [[2887. 58. 63.]  
 [ 55. 2932. 8.]  
 [ 58. 10. 3929.]]



The decision boundaries appear to be smooth and curved surfaces

## PART B.

### ERM Classification with Loss Matrices

2 loss matrices  $\Lambda_{10}$  and  $\Lambda_{100}$  are defined to penalize misclassifications more heavily when  $L = 3$ .

$$\Lambda_{10} = \begin{bmatrix} 0 & 10 & 10 \\ 1 & 0 & 10 \\ 1 & 1 & 0 \end{bmatrix}; \quad \Lambda_{100} = \begin{bmatrix} 0 & 100 & 100 \\ 1 & 0 & 100 \\ 1 & 1 & 0 \end{bmatrix}$$

### Expected Risk minimization

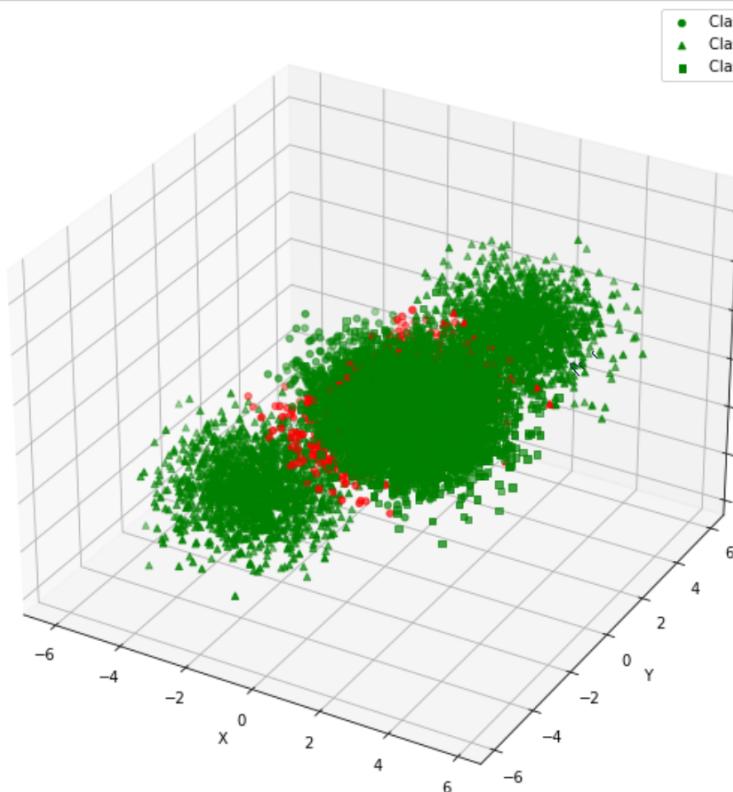
The classification rule that minimizes the expected risk for each sample:

$$\hat{L}(x) = \arg \min_i \sum_{j=1}^3 \Lambda(i, j) P(L=j|x)$$

This means that for each class  $i$ , we compute the expected risk by summing the posterior probabilities  $P(L=j|x)$  weighted by the corresponding loss  $\Lambda(i, j)$  and choose the class  $i$  that minimizes the total risk.

#### Part B Results:

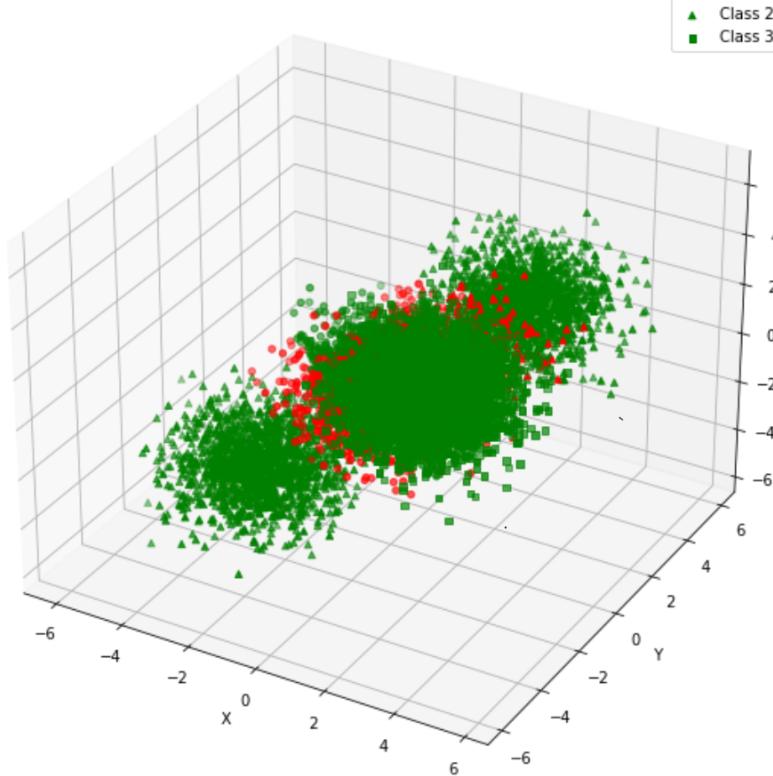
ERM Classifier (Loss Matrix 10) Accuracy: 0.9565  
 ERM Classifier (Loss Matrix 100) Accuracy: 0.8883



- Class 1
- ▲ Class 2
- Class 3

Introduces moderate bias, slightly favoring class 3.

z Higher cost for misclassifying class 3.



Dramatic shift in the decision boundary.

Class 1 has reduced significantly

Class 3 dominates a larger portion of the feature space.

High penalties for misclassifying class 1. Misclassification into class 3 are less costly.

The Bayes classifier achieves a better overall performance with balanced error rates across classes. The ERM classifier shows a trade off between accuracy and cost sensitivity. With loss matrix 10 there's a minor (1.83%) accuracy sacrifice for moderate cost optimization. With Loss matrix there's a significant accuracy drop suggesting strong avoidance of costly errors.

The progression from Bayes to ERM with increasing loss values demonstrates how the decision boundaries can be manipulated to address specific cost structure in classification problems.

Q3.

To implement minimum probability of error classifiers for the given Wine Quality and Human Activity Recognition datasets we follow these steps:

1. The data is preprocessed  $\rightarrow$  load and clean the data the split the data into features and labels.
2. Estimate mean vectors and covariance matrices for each class.

Apply regularization to covariance matrices by adding a small value  $\lambda$  to the diagonal if the covariance matrix is ill-conditioned (i.e. has very small or zero eigenvalues)

$$\Sigma_{\text{regularized}} = \Sigma_{\text{sample}} + \lambda I$$

3. Estimate Class priors by counting the number of samples per class
4. For each class, assume that the features follow a multivariate gaussian distribution

$$P(x|y) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_y|}} \exp\left(-\frac{1}{2} (x - \mu_y)^T \Sigma_y^{-1} (x - \mu_y)\right)$$

5. Using the gaussian model for each class and the class priors classify the samples based on maximum posterior probability

$$P(y|x) = P(x|y) P(y)$$

6. Evaluate the error probability and confusion matrix and visualize using principal components (to reduce the dimensionality of data)

$$P(\text{error}) = \frac{\text{number of misclassifications}}{\text{total number of samples}}$$

To evaluate the performance of the classifier.

We assume that different classes are independent of each other, allowing us to estimate parameters separately for each class. The model doesn't assume independence between features with a class as it uses full covariance matrix.

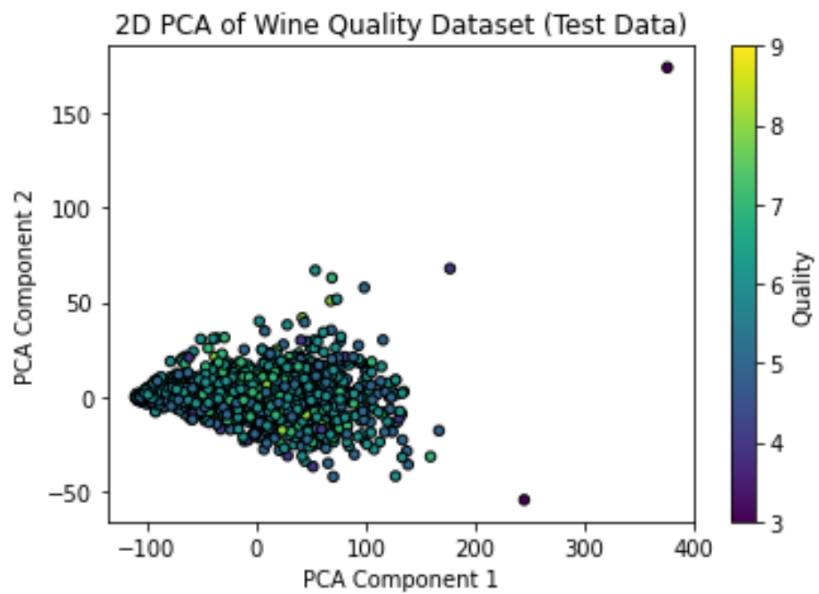
# WINE

```
In [10]: M 1 import numpy as np
2 import pandas as pd
3 from scipy.stats import multivariate_normal
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import confusion_matrix, accuracy_score
6 import matplotlib.pyplot as plt
7 from sklearn.decomposition import PCA
8
9 # Load the dataset (white wine in this case, but it can be easily changed to red wine)
10 white_wine_data = pd.read_csv("C:/Users/aksha/Downloads/wine+quality/winequality-red.csv", sep=";")
11 red_wine_data = pd.read_csv("C:/Users/aksha/Downloads/wine+quality/winequality-white.csv", sep=";")
12
13 # Combine the two datasets into one (if needed)
14 wine_data = pd.concat([white_wine_data, red_wine_data])
15
16 # Separate features (X) and Labels (y)
17 X = wine_data.iloc[:, :-1].values # All columns except the last one (quality)
18 y = wine_data.iloc[:, -1].values # The last column (quality)
19
20 # Split the data into training and testing sets (80% training, 20% testing)
21 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
22
23 # Estimate the mean vectors and covariance matrices for each class
24 class_labels = np.unique(y_train) # Unique class labels (wine quality scores)
25 means = {}
26 covariances = {}
27 priors = {}
28
29 # Regularization parameter (based on sample covariance matrix eigenvalues)
30 def regularize_covariance(cov_matrix, alpha=0.01):
31     trace = np.trace(cov_matrix)
32     rank = np.linalg.matrix_rank(cov_matrix)
33     lambda_reg = alpha * trace / rank
34     return cov_matrix + lambda_reg * np.eye(cov_matrix.shape[0])
35
36 for label in class_labels:
37     X_class = X_train[y_train == label]
38     means[label] = np.mean(X_class, axis=0)
39     covariances[label] = np.cov(X_class, rowvar=False)
40
41 # Regularize the covariance matrix if necessary
42 covariances = regularize_covariance(covariances)
43
44 # Estimate class prior probabilities
45 priors[label] = len(X_class) / len(X_train)
46
47 # Minimum-P(error) classification rule
48 def classify(sample, means, covariances, priors):
49     posteriors = []
50     for label in class_labels:
51         likelihood = multivariate_normal.pdf(sample, mean=means[label], cov=covariances[label])
52         posterior = likelihood * priors[label]
53         posteriors.append(posterior)
54
55     return class_labels[np.argmax(posteriors)]
56
57 # Apply the classifier to the test data
58 y_pred = [classify(sample, means, covariances, priors) for sample in X_test]
59
60 # Compute the confusion matrix and error probability
61 conf_matrix = confusion_matrix(y_test, y_pred)
62 error_probability = 1 - accuracy_score(y_test, y_pred)
63
64 print("Confusion Matrix:")
65 print(conf_matrix)
66 print(f"Error Probability: {error_probability:.4f}")
67
68 # Visualize the dataset using PCA (2D and 3D projections)
69 def plot_pca(X, y, title, n_components=2):
70     pca = PCA(n_components=n_components)
71     X_pca = pca.fit_transform(X)
72
73     if n_components == 2:
74         plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y, cmap='viridis', edgecolor='k', s=20)
75         plt.title(title)
76         plt.xlabel('PCA Component 1')
77         plt.ylabel('PCA Component 2')
78         plt.colorbar(label='Quality')
79         plt.show()
80     elif n_components == 3:
81         fig = plt.figure()
82         ax = fig.add_subplot(111, projection='3d')
83         ax.scatter(X_pca[:, 0], X_pca[:, 1], X_pca[:, 2], c=y, cmap='viridis', edgecolor='k', s=20)
84         ax.set_title(title)
85         ax.set_xlabel('PCA Component 1')
86         ax.set_ylabel('PCA Component 2')
87         ax.set_zlabel('PCA Component 3')
88         plt.show()
89
90 # Plot the dataset using 2D and 3D PCA projections
91 plot_pca(X_test, y_test, '2D PCA of Wine Quality Dataset (Test Data)', n_components=2)
92 plot_pca(X_test, y_test, '3D PCA of Wine Quality Dataset (Test Data)', n_components=3)
```

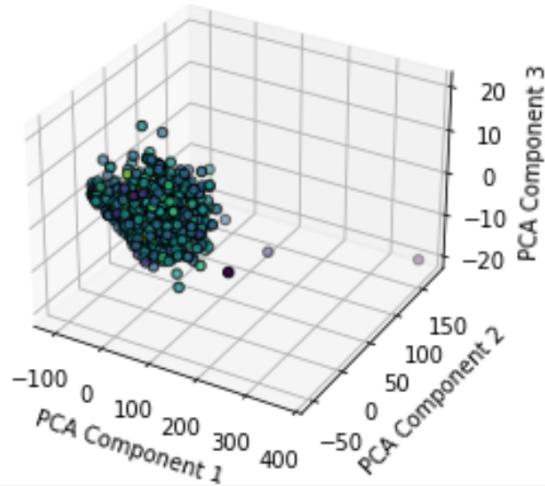
Confusion Matrix:

```
[[ 1  1  0  3  1  0  0]
 [ 1  0  2  29 11  0  0]
 [ 1  1  6 300 93  1  0]
 [ 1  1  3 316 275  1  0]
 [ 1  0  1  58 155  0  0]
 [ 0  0  0  9  25  2  0]
 [ 0  0  0  0  1  0  0]]
```

Error Probability: 0.6308



3D PCA of Wine Quality Dataset (Test Data)



### Interpretation:

- The error probability of 63.08% is quite high indicating that the gaussian classifier is not performing well on this dataset. This suggests that the assumption of gaussian distribution for wine quality classes may not be appropriate.
- From the confusion matrix we see a strong diagonal pattern in some rows indicating better prediction for these classes. However there's significant off-diagonal confusion. For instance row 3 has 300 correct predictions but confuses 93 samples with the next class.

Similarly class 5 has 316 correct predictions but confuses 275 samples with the next class. Some classes like the last row have very few samples leading to poor prediction.

From the 2D PCA plot we see significant overlap between different wine quality classes. The data doesn't form clear, separate clusters. There's a gradual color transition, indicating that the wine quality is more of a continuum classes.

The 3D PCA plot shows that even with an additional dimension, classes remain heavily overlapped.

Conclusion:

The gaussian assumption is likely not ideal for this dataset. Wine quality appears to be more of a continuous spectrum rather than distinct classes. The high error rate suggests that using discrete classification might not be the best approach.

We could treat this as a regression problem instead of classification by using ordinal classification techniques that account for the ordered nature of wine quality scores to get better results. The regularization term ( $\lambda I$ ) could be tuned further for better results.

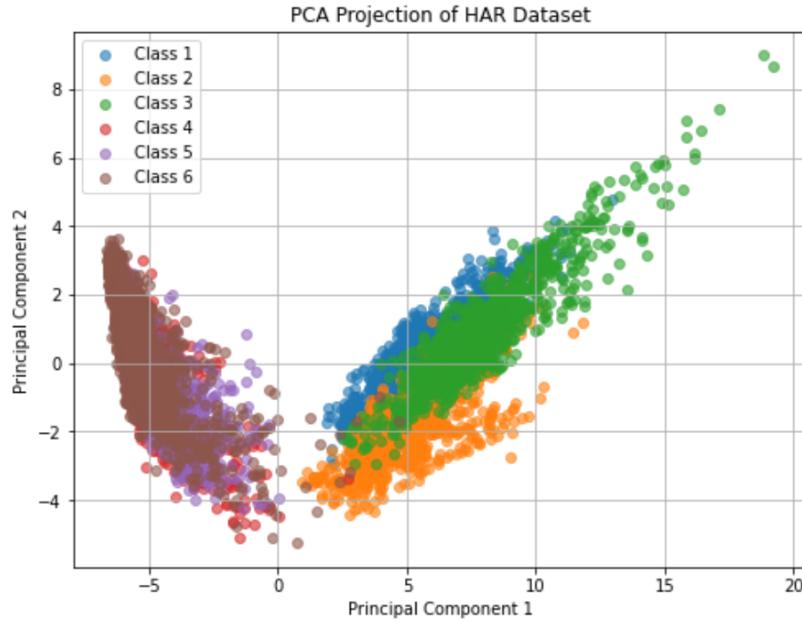
# Human Activity Recognition

```
In [17]: M
 1 import pandas as pd
 2 import numpy as np
 3
 4 # Load features and labels
 5 X_train = np.loadtxt('C:/Users/aksha/Downloads/UCI HAR Dataset/train/X_train.txt')
 6 y_train = np.loadtxt('C:/Users/aksha/Downloads/UCI HAR Dataset/train/y_train.txt').astype(int)
 7 X_test = np.loadtxt('C:/Users/aksha/Downloads/UCI HAR Dataset/test/X_test.txt')
 8 y_test = np.loadtxt('C:/Users/aksha/Downloads/UCI HAR Dataset/test/y_test.txt').astype(int)
 9
10 # Check the shape of the data
11 print("Training data shape:", X_train.shape)
12 print("Training labels shape:", y_train.shape)
13 print("Test data shape:", X_test.shape)
14 print("Test labels shape:", y_test.shape)
15
16
17 def estimate_mean_cov(X, y, num_classes, regularization_lambda=0.01):
18     means = []
19     covariances = []
20     priors = []
21
22     for c in range(1, num_classes + 1):
23         X_c = X[y == c] # Data belonging to class c
24         mean_c = np.mean(X_c, axis=0)
25         cov_c = np.cov(X_c, rowvar=False)
26
27         # Apply regularization to the covariance matrix
28         cov_c_reg = cov_c + regularization_lambda * np.eye(cov_c.shape[0])
29
30         means.append(mean_c)
31         covariances.append(cov_c_reg)
32         priors.append(len(X_c) / len(y)) # Prior probability of class c
33
34     return np.array(means), np.array(covariances), np.array(priors)
35
36 # Estimate mean vectors, covariance matrices, and priors
37 num_classes = 6 # There are 6 activity labels
38 means, covs, priors = estimate_mean_cov(X_train, y_train, num_classes)
39
40 # Check shapes
41 print("Means shape:", means.shape)
42 print("Covariances shape:", covs.shape)
43 print("Priors:", priors)
44
45
46 from scipy.stats import multivariate_normal
47
48 def classify_gaussian(X, means, covs, priors):
49     num_samples = X.shape[0]
50     num_classes = len(means)
51     log_probs = np.zeros((num_samples, num_classes))
52
53     for c in range(num_classes):
54         mvn = multivariate_normal(mean=means[c], cov=covs[c])
55         log_probs[:, c] = mvn.logpdf(X) + np.log(priors[c])
56
57     return np.argmax(log_probs, axis=1) + 1 # Class with highest posterior probability
58
59 # Classify the test data
60 y_pred = classify_gaussian(X_test, means, covs, priors)
61
62
63 from sklearn.metrics import confusion_matrix, accuracy_score
64
65 # Confusion matrix
66 conf_matrix = confusion_matrix(y_test, y_pred)
67 print("Confusion Matrix:\n", conf_matrix)
68
69 # Accuracy score
70 accuracy = accuracy_score(y_test, y_pred)
71 print("Accuracy:", accuracy)
72
73 from sklearn.decomposition import PCA
74 import matplotlib.pyplot as plt
75
76 # Reduce dimensionality to 2 for visualization
77 pca = PCA(n_components=2)
78 X_train_pca = pca.fit_transform(X_train)
79
80 # Plot the projected data
81 plt.figure(figsize=(8, 6))
82 for label in np.unique(y_train):
83     plt.scatter(X_train_pca[y_train == label, 0], X_train_pca[y_train == label, 1], label=f'Class {label}', alpha=0.6)
84
85 plt.title('PCA Projection of HAR Dataset')
86 plt.xlabel('Principal Component 1')
87 plt.ylabel('Principal Component 2')
88 plt.legend()
89 plt.grid(True)
90 plt.show()
```

```

Training data shape: (7352, 561)
Training labels shape: (7352,)
Test data shape: (2947, 561)
Test labels shape: (2947,)
Means shape: (6, 561)
Covariances shape: (6, 561, 561)
Priors: [0.16675734 0.14594668 0.13411317 0.17491839 0.18688792 0.1913765 ]
Confusion Matrix:
[[482   1  13   0   0   0]
 [ 0 470   1   0   0   0]
 [ 4  41 375   0   0   0]
 [ 0   0   0 387 103   1]
 [ 0   0   0   7 525   0]
 [ 0   0   0   0  0 537]]
Accuracy: 0.9419748897183576

```



Strong diagonal elements indicate good classification for most classes.

41 samples from Class 3 misclassified as class 2

103 samples from class 4 misclassified as class 5.

The PCA projection plot shows good separation between some classes. The model achieved an accuracy of approximately 94.2%. This suggests that the gaussian model works well for most classes. The separation in PCA space for some classes aligns with gaussian assumption.

High dimensionality (561 features) may impact covariance estimation. Dimension reduction techniques might be beneficial for computational efficiency & improved classification by focusing on the most deterministic features.

The use of  $\lambda I$  regularization was crucial for numerical stability. Priors seem relatively balanced [0.167, 0.146, 0.134, 0.175, 0.187, 0.191]

Codes and Results :

<https://colab.research.google.com/drive/1u88HJsaOSaJFS3n06uqukpz8fLsYBj5M?usp=sharing>