

# INTRODUCTION TO MACHINE LEARNING

AKSHATA KUMBLE

Q1.

We will create a multi-class classification problem with 4 classes and gaussian distributions. Using a theoretically optimal classifier as a benchmark we will train multiple MLPs with a variety of configurations, using cross-entropy loss and maximum likelihood estimation. This allows us to approximate the MAP (Maximum a Posteriori) classification rule.

Process of developing the solution.

1. Define data distribution and generate samples

1. Define Class Conditional Distribution:

We have  $C=4$  classes and for each class we define a 3 dimensional gaussian distribution with a unique mean vector and covariance matrix. Since we want the theoretically optimal classifier to have a probability of error between 10-20%, we set mean vectors and covariance matrices such that the classes overlap slightly, that should lead to some classification ambiguity.

2. Generating Training and Test sets

We generate several datasets for training (with various sample sizes) and a large test dataset for performance evaluation. For each class, we generate samples from the gaussian distribution defined for that class. Combine the samples from all classes and assign class labels. Generate training sets with 100, 500, 5000, 10000 samples and test set with 100,000 samples.

To test if this setup gives us desired error rates for a MAP classifier, we can generate a large dataset and apply a MAP classifier with knowledge of true parameters. The MAP classifier will compute the likelihood of each class & assign the sample to the class with the highest posterior probability.

With the specified parameters, we should observe an error rate b/w 10-20% depending on the degree of overlap in our gaussian distribution.

If the error rate is outside this range we can iteratively adjust the means or covariances by bringing the mean vectors closer or increasing the off diagonal elements in the covariance matrices for added correlation.

Adjustments:- ① Means: chose positions relatively close to each other to create class overlap. This overlap increases misclassification achieving the desired error rate. ② Covariances: Used varying covariance sizes to control the spread, introducing slight asymmetry with off diag elements makes certain regions in feature space ambiguous. Expanding the variances will increase the spread of each class making them overlap more. We can increase the diagonal values making the covariance matrices more asymmetrical.

```

1 import numpy as np
2
3 # Adjusted mean vectors, not too close or too far
4 class_means = [
5     np.array([1.5, 1.5, 1.5]),
6     np.array([-1.5, -1.5, -1.5]),
7     np.array([1.5, -1.5, 1.5]),
8     np.array([-1.5, 1.5, -1.5])
9 ]
10
11 # Adjusted covariance matrices with moderate variances
12 class_covariances = [
13     np.array([[1.8, 0.4, 0.2], [0.4, 1.8, 0.3], [0.2, 0.3, 1.8]]),
14     np.array([[1.9, 0.3, 0.1], [0.3, 1.9, 0.2], [0.1, 0.2, 1.9]]),
15     np.array([[1.8, -0.3, 0.5], [-0.3, 1.7, 0.2], [0.5, 0.2, 1.6]]),
16     np.array([[1.7, 0.5, -0.3], [0.5, 1.7, 0.6], [-0.3, 0.6, 1.8]])
17 ]
18
19 from scipy.stats import multivariate_normal
20
21 # Generate test dataset to evaluate error probability
22 def generate_data(n_samples, class_means, class_covariances):
23     X = []
24     y = []
25     for i, (mean, cov) in enumerate(zip(class_means, class_covariances)):
26         X_class = np.random.multivariate_normal(mean, cov, n_samples // len(class_means))
27         X.append(X_class)
28         y.append(np.full(X_class.shape[0], i))
29     return np.vstack(X), np.hstack(y)
30
31 # Create a test dataset of 100,000 samples
32 test_data, test_labels = generate_data(100000, class_means, class_covariances)
33
34
35 # MAP classifier function using known class pdfs
36 def theoretical_MAP_classifier(X, class_means, class_covariances):
37     posteriors = []
38     for mean, cov in zip(class_means, class_covariances):
39         # Calculate posterior probability assuming equal priors
40         posterior = multivariate_normal(mean=mean, cov=cov).pdf(X)
41         posteriors.append(posterior)
42     # Assign each sample to the class with the highest posterior probability
43     return np.argmax(posteriors, axis=0)
44
45 # Get predictions and calculate probability of error
46 theoretical_preds = theoretical_MAP_classifier(test_data, class_means, class_covariances)
47 theoretical_error = np.mean(theoretical_preds != test_labels)
48 print(f"Theoretical MAP Classifier Probability of Error: {theoretical_error:.2%}")
49

```

Theoretical MAP Classifier Probability of Error: 17.08%

```

1 import numpy as np
2
3 class_means = [
4     np.array([1.5, 1.5, 1.5]),
5     np.array([-1.5, -1.5, -1.5]),
6     np.array([1.5, -1.5, 1.5]),
7     np.array([-1.5, 1.5, -1.5])
8 ]
9
10 # Adjusted covariance matrices with moderate variances
11 class_covariances = [
12     np.array([[1.8, 0.4, 0.2], [0.4, 1.8, 0.3], [0.2, 0.3, 1.8]]),
13     np.array([[1.9, 0.3, 0.1], [0.3, 1.9, 0.2], [0.1, 0.2, 1.9]]),
14     np.array([[1.8, -0.3, 0.5], [-0.3, 1.7, 0.2], [0.5, 0.2, 1.6]]),
15     np.array([[1.7, 0.5, -0.3], [0.5, 1.7, 0.6], [-0.3, 0.6, 1.8]])
16 ]
17
18 # Set random seed for reproducibility
19 np.random.seed(42)
20
21 # Function to generate samples for a single dataset size
22 def generate_dataset(n_samples):
23     X = []
24     y = []
25
26     samples_per_class = n_samples // 4 # Uniform prior assumption (each class has equal samples)
27
28     for class_index, (mean, cov) in enumerate(zip(class_means, class_covariances)):
29         # Sample from Gaussian distribution for this class
30         samples = np.random.multivariate_normal(mean, cov, samples_per_class)
31         labels = np.full(samples_per_class, class_index) # Label all samples with the current class
32
33         X.append(samples)
34         y.append(labels)
35
36     # Concatenate all class samples to form the full dataset
37     X = np.vstack(X)
38     y = np.concatenate(y)
39
40     # Shuffle the dataset (helps to randomize order)
41     indices = np.random.permutation(n_samples)
42     X, y = X[indices], y[indices]
43
44     return X, y
45
46 # Generate datasets of specified sizes
47 dataset_sizes = [100, 500, 1000, 5000, 10000]
48 train_datasets = {size: generate_dataset(size) for size in dataset_sizes}
49
50 # Generate a large test dataset of 100000 samples for evaluation
51 test_X, test_y = generate_dataset(100000)
52
53 # Show sample sizes and a preview
54 for size, (X, y) in train_datasets.items():
55     print(f"Training dataset with {size} samples: X shape = {X.shape}, y shape = {y.shape}")
56
57 print(f"Test dataset: X shape = {test_X.shape}, y shape = {test_y.shape}")
58

```

Training dataset with 100 samples: X shape = (100, 3), y shape = (100,)  
 Training dataset with 500 samples: X shape = (500, 3), y shape = (500,)  
 Training dataset with 1000 samples: X shape = (1000, 3), y shape = (1000,)  
 Training dataset with 5000 samples: X shape = (5000, 3), y shape = (5000,)  
 Training dataset with 10000 samples: X shape = (10000, 3), y shape = (10000,)  
 Test dataset: X shape = (100000, 3), y shape = (100000,)

## 2. Theoretically Optimal Classifier (MAP classifier)

Using the known class conditional distributions, we can construct a theoretical MAP classifier.

### Compute Posterior Probabilities

Since the priors are uniform, the posterior probability for each class  $C_i$  is proportional to likelihood  $P(x|C_i)$ . For each sample  $x$ , compute the likelihoods for all classes using the gaussian PDFs, then normalize them to obtain posterior probabilities.

### Decision Rule

Assign each sample  $x$  to the class with the highest posterior probability. Apply this rule to classify the test datasets and calculate the empirical probability of error for this optimal classifier, which will act as a baseline for evaluating our MLP models.

To construct the minimum probability of error classification rule using the true pdf we need to use the maximum a posteriori rule. The MAP classification rule chooses the class  $C_k$  with the highest posterior probability for a given sample  $x$ :

$$C_k = \arg \max_k P(C_k|x) = \arg \max_k P(x|C_k) P(C_k)$$

since the priors are uniform (ie  $P(C_k) = \frac{1}{4}$  for all  $k$ ) the MAP rule simplifies to:

$$C_k = \arg \max_k P(x|C_k) \quad \text{where } P(x|C_k) \text{ is the class conditional}$$

PDF for class  $C_k$  which is gaussian with mean vector  $\mu_k$  & covariance matrix  $\Sigma_k$  of dimensionality  $n$ .

Class Conditional distributions :  $P(x|C_k) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)\right)$  where  $d=3$

Apply MAP rule for each sample in the test dataset & select class  $C_k$  that maximizes the likelihood.

Compare the predicted labels with the true labels and calculate the empirical probability of error.

```
1 import numpy as np
2 from scipy.stats import multivariate_normal
3
4 # Function to calculate the likelihood of x given class k (Gaussian pdf)
5 def calculate_likelihood(x, class_index):
6     mean = class_means[class_index]
7     cov = class_covariances[class_index]
8
9     # Multivariate normal pdf
10    return multivariate_normal.pdf(x, mean=mean, cov=cov)
11
12 # MAP classifier: Choose the class with the highest likelihood
13 def map_classifier(x, class_means, class_covariances):
14     max_likelihood = -np.inf
15     predicted_class = -1
16     for class_index in range(4):
17         likelihood = calculate_likelihood(x, class_index)
18         if likelihood > max_likelihood:
19             max_likelihood = likelihood
20             predicted_class = class_index
21     return predicted_class
22
23 # Classify test set and calculate error rate
24 def calculate_map_error(test_X, test_y, class_means, class_covariances):
25     correct_predictions = 0
26     for x, true_label in zip(test_X, test_y):
27         predicted_label = map_classifier(x, class_means, class_covariances)
28         if predicted_label == true_label:
29             correct_predictions += 1
30
31     # Probability of error is the proportion of incorrect predictions
32     error_rate = 1 - (correct_predictions / len(test_y))
33     return error_rate
34
35 # Apply MAP classifier to the test set
36 map_error = calculate_map_error(test_X, test_y, class_means, class_covariances)
37 print(f"Theoretical MAP Classifier Probability of Error: {map_error * 100:.2f}%")
```

(16.93%)

This error rate represents the aspirational performance level for the MLP classifier.

### 3. Multilayer Perceptron (MLP) Structure

We need a 2 layer MLP structure : 1. Hidden Layer : A single hidden layer with P perceptions where P is determined using cross validation. Use smooth ramp activation function like ReLU to add non-linearity. 2. Output Layer : Use a softmax function at the output layer to produce probabilities for each class. The input layer accepts a 3D input vector  $x$ .

Input Transformation :  $h^{(l)} = f(W^{(l)}h^{(l-1)} + b^{(l)})$  where  $h^{(l)}$  is the i-th layer,  $W^{(l)}$  &  $b^{(l)}$  are the weight matrix and bias vector for layer l,  $f$  is a non-linear activation function (ReLU). Output layer applies softmax function to convert the scores into probabilities for each class.

$$P(y=c|x) = \frac{e^{z_c}}{\sum_j e^{z_j}}$$
 where  $z_c$  is the score for class  $c$  in the output layer.

The next steps are to train the MLP models, use cross validation to tune the number of hidden perceptions for optimal performance

```
1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 from sklearn.model_selection import KFold
5 from sklearn.metrics import accuracy_score
6
7 # Define the MLP model with ReLU activation in the hidden Layer
8 class MLP(nn.Module):
9     def __init__(self, input_size, hidden_size, output_size):
10         super(MLP, self).__init__()
11         self.hidden = nn.Linear(input_size, hidden_size)
12         self.activation = nn.ReLU()
13         self.output = nn.Linear(hidden_size, output_size)
14         self.softmax = nn.Softmax(dim=1)
15
16     def forward(self, x):
17         x = self.activation(self.hidden(x))
18         x = self.softmax(self.output(x))
19         return x
20
21 # Hyperparameters and configurations
22 input_size = 3 # dimension of the input feature vector
23 output_size = 4 # Number of classes
24 hidden_sizes = [5, 10, 20, 50] # candidate sizes for cross-validation
25 num_folds = 10 # Number of cross-validation folds
26
27 # Loss function (Cross-Entropy Loss) and function to train the MLP
28 loss_fn = nn.CrossEntropyLoss()
```

### 4. Model order Selection and Training.

#### Model Selection via Cross-Validation :

For each training set, we perform 10 fold cross validation to identify the optimal number of perceptions in the hidden layer that minimizes classification error probability. We use cross validation for hyper-tuning the MLP structure. To perform 10 fold cross validation to select the best number of perceptions for the hidden layer of the MLP, we split the training data into 10 equal parts and use 9 parts for training while the remaining part will be used for validation. This will be repeated 10 times. The objective is to minimize the classification error. For each number of perceptions in the hidden layer, we will calculate the average classification error over 10 folds. We then repeat the process for different numbers of perceptions in the hidden layer and select the value that yields the lowest cross validation error.

Training MLPs : Using the cross validated number of perceptions, train the MLP on each training set. Use maximum likelihood parameter estimation and randomly reinitialize the weights multiple times to avoid local minima, retaining the model with the best log-likelihood.

We used KFold from sklearn.model\_selection for a 10 fold cross validation. For each fold, we trained an MLP model with different hidden layer sizes (from 1 to 20 perceptions) and evaluate their error rates. StandardScaler was used to normalize the data before training the MLPs. Our model uses ReLU activation function and Adam optimizer.

Once the optimal number of perceptions was identified via cross validation for each dataset size, we train an MLP on the entire dataset. We performed the training for each dataset size using the best performing number of perceptions and evaluated

the test set error rate by comparing the predicted labels with the true labels. For each trained MLP we estimated the empirical

probability of error by predicting the test set labels & calculating the misclassification rate. Error rate = 1 - accuracy\_score

```
1 import numpy as np
2 from sklearn.model_selection import KFold
3 from sklearn.neural_network import MLPClassifier
4 from sklearn.metrics import accuracy_score
5 from sklearn.preprocessing import StandardScaler
6
7 # Step 1: Dataset generation function
8 class_means = [
9     np.array([1.5, 1.5, 1.5]),
10    np.array([-1.5, -1.5, -1.5]),
11    np.array([1.5, -1.5, 1.5]),
12    np.array([-1.5, 1.5, -1.5])
13 ]
14
15 # Adjusted covariance matrices with moderate variances
16 class_covariances = [
17     np.array([[1.8, 0.4, 0.2], [0.4, 1.8, 0.3], [0.2, 0.3, 1.8]]),
18     np.array([[1.9, 0.3, 0.1], [0.3, 1.9, 0.2], [0.1, 0.2, 1.9]]),
19     np.array([[1.8, -0.3, 0.5], [-0.3, 1.7, 0.2], [0.5, 0.2, 1.6]]),
20     np.array([[1.7, 0.5, -0.3], [0.5, 1.7, 0.6], [-0.3, 0.6, 1.8]])
21 ]
22
23 # Set random seed for reproducibility
24 np.random.seed(42)
25
26 # Function to generate samples for a single dataset size
27 def generate_dataset(n_samples):
28     X = []
29     y = []
30
31     samples_per_class = n_samples // 4 # Uniform prior assumption (each class has equal samples)
32
33     for class_index, (mean, cov) in enumerate(zip(class_means, class_covariances)):
34         # Sample from Gaussian distribution for this class
35         samples = np.random.multivariate_normal(mean, cov, samples_per_class)
36         labels = np.full(samples_per_class, class_index) # Label all samples with the current class
37
38         X.append(samples)
39         y.append(labels)
40
41     # Concatenate all class samples to form the full dataset
42     X = np.vstack(X)
43     y = np.concatenate(y)
44
45     # Shuffle the dataset (helps to randomize order)
46     indices = np.random.permutation(n_samples)
47     X, y = X[indices], y[indices]
48
49     return X, y
50
51 # Generate datasets of specified sizes
52 dataset_sizes = [100, 500, 1000, 5000, 10000]
53 train_datasets = {size: generate_dataset(size) for size in dataset_sizes}
54
55 # Generate a large test dataset of 100000 samples for evaluation
56 test_X, test_y = generate_dataset(100000)
57
58 # Show sample sizes and a preview
59 for size, (X, y) in train_datasets.items():
60     print(f"Training dataset with {size} samples: X shape = {X.shape}, y shape = {y.shape}")
61
62 print(f"Test dataset: X shape = {test_X.shape}, y shape = {test_y.shape}")
63
64 # Step 2: Cross-validation function for MLP
65 def cross_validate_mlp(X, y, max_perceptrons, n_splits=10):
66     kf = KFold(n_splits=n_splits, shuffle=True, random_state=42)
67     best_perceptrons = None
68     best_error = float('inf')
69
70     # List to store errors for each perceptron count
71     error_rates = []
72
73     for P in range(1, max_perceptrons + 1):
74         fold_errors = []
75
76         for train_idx, val_idx in kf.split(X):
77             X_train, X_val = X[train_idx], X[val_idx]
78             y_train, y_val = y[train_idx], y[val_idx]
79
80             # Standardize the data (important for MLPs)
81             scaler = StandardScaler()
82             X_train_scaled = scaler.fit_transform(X_train)
83             X_val_scaled = scaler.transform(X_val)
84
85             # Initialize MLP with current number of perceptrons in the hidden layer
86             mlp = MLPClassifier(hidden_layer_sizes=(P,), max_iter=1000, random_state=42, solver='adam', activation='relu')
87             mlp.fit(X_train_scaled, y_train)
88
89             # Make predictions on the validation set
90             y_pred = mlp.predict(X_val_scaled)
91             fold_error = 1 - accuracy_score(y_val, y_pred) # Probability of error
92             fold_errors.append(fold_error)
93
94             # Calculate average error over all folds for the current P
95             avg_error = np.mean(fold_errors)
96             error_rates.append(avg_error)
97
98             # Update best perceptrons if the current error is lower
99             if avg_error < best_error:
100                 best_error = avg_error
101                 best_perceptrons = P
102
103     return best_perceptrons, error_rates
104
105 # Step 3: Example usage and cross-validation for each dataset
106 max_perceptrons = 20 # Set the maximum number of perceptrons to test
107 best_perceptrons = []
108
109 for size, (X_train, y_train) in train_datasets.items():
110     print(f"Performing cross-validation for dataset with {size} samples...")
111     best_P, error_rates = cross_validate_mlp(X_train, y_train, max_perceptrons)
112     best_perceptrons[size] = best_P
113     print(f"\nBest number of perceptrons for {size} samples: {best_P}\n")
114     print("Error rates for different perceptrons:\n", error_rates)
115
116 # Step 4: Train MLP using the best number of perceptrons
117 # Using the training set with the best number of perceptrons
118 for size, (X_train, y_train) in train_datasets.items():
119     best_P = best_perceptrons[size]
120     print(f"\nTraining MLP with {best_P} perceptrons for {size} samples...\n")
121
122     # Standardize the entire training set and train the model
123     scaler = StandardScaler()
124     X_train_scaled = scaler.fit_transform(X_train)
125
126     # Train the MLP with the optimal number of perceptrons
127     best_mlp = MLPClassifier(hidden_layer_sizes=(best_P,), max_iter=1000, random_state=42, solver='adam', activation='relu')
128     best_mlp.fit(X_train_scaled, y_train)
129
130     # Step 5: Evaluate performance on test dataset
131     print("Evaluating MLP on test dataset...")
132     X_test_scaled = scaler.transform(test_X)
133     y_pred = best_mlp.predict(X_test_scaled)
134     test_error = 1 - accuracy_score(test_y, y_pred)
135     print(f"\nTest Error Rate: {test_error:.4f}\n")
```

Training MLP with 13 perceptrons for 100 samples...
Evaluating MLP on test dataset...
Test Error Rate: 0.1921

Training MLP with 13 perceptrons for 500 samples...
Evaluating MLP on test dataset...
Test Error Rate: 0.1779

Training MLP with 5 perceptrons for 1000 samples...
Evaluating MLP on test dataset...
Test Error Rate: 0.1771

Training MLP with 20 perceptrons for 5000 samples...
Evaluating MLP on test dataset...
Test Error Rate: 0.1718

Training MLP with 9 perceptrons for 10000 samples...
Evaluating MLP on test dataset...
Test Error Rate: 0.1711

Performing cross-validation for dataset with 100 samples...

Best number of perceptrons for 100 samples: 13

Error rates for different perceptrons: [0.56, 0.2400000000000005, 0.1499999999999997, 0.18, 0.1699999999999998, 0.1699999999999998, 0.1699999999999998, 0.1599999999999998, 0.1500000000000002, 0.1599999999999998, 0.18, 0.1699999999999998, 0.14, 0.1499999999999997, 0.14, 0.1699999999999998, 0.1500000000000002, 0.1599999999999998, 0.16, 0.1599999999999998]

Performing cross-validation for dataset with 500 samples...

Best number of perceptrons for 500 samples: 13

Error rates for different perceptrons: [0.438, 0.17, 0.158, 0.164, 0.1559999999999997, 0.152, 0.1599999999999998, 0.164, 0.16, 0.158, 0.156, 0.162, 0.15, 0.1500000000000005, 0.162, 0.156, 0.16, 0.16, 0.156, 0.16]

Performing cross-validation for dataset with 1000 samples...

Best number of perceptrons for 1000 samples: 5

Error rates for different perceptrons: [0.3919999999999996, 0.1750000000000002, 0.175, 0.18, 0.17, 0.1730000000000001, 0.1750000000000002, 0.1780000000000002, 0.175, 0.1740000000000002, 0.1720000000000001, 0.171, 0.178, 0.1740000000000002, 0.1709999999999999, 0.1770000000000002, 0.1790000000000002, 0.1750000000000002, 0.171, 0.1800000000000002]

Performing cross-validation for dataset with 5000 samples...

Best number of perceptrons for 5000 samples: 20

Error rates for different perceptrons: [0.3706, 0.1686, 0.1668, 0.1689999999999998, 0.1692000000000002, 0.1666000000000003, 0.1662000000000004, 0.1672000000000007, 0.1710000000000004, 0.1718000000000004, 0.1704000000000002, 0.1670000000000004, 0.1668, 0.1668, 0.1670000000000004, 0.1694000000000002, 0.1664000000000002, 0.1684000000000002, 0.1680000000000007, 0.1656000000000002]

Performing cross-validation for dataset with 10000 samples...

Best number of perceptrons for 10000 samples: 9

Error rates for different perceptrons: [0.3469000000000004, 0.1696000000000003, 0.1699000000000002, 0.1659000000000002, 0.1654000000000005, 0.1653, 0.1667000000000001, 0.1662000000000004, 0.1625, 0.1639000000000002, 0.165, 0.1634000000000002, 0.1643000000000003, 0.1633000000000006, 0.1655000000000004, 0.1634000000000002, 0.1629000000000002, 0.1644000000000002, 0.1633000000000006, 0.1642000000000004]

These results show the outcome of running cross validation to determine the best no. of perceptrons for the hidden layer of the MLP.

We observe that as the number of samples increases, the optimal number of perceptron varies, reflecting the complexity of the dataset. Smaller datasets (like 100 and 500 samples) are likely to benefit from more perceptrons because they help the model capture more detailed patterns, while larger datasets (like 5000 and 10000 samples) may need fewer perceptrons, which helps avoid overfitting and ensures better generalization.

The error rates for different number of perceptrons represents how well the MLP performs for various configurations. For instance, in the 100 samples dataset the error starts at 0.56 for 1 perceptron and decreases to 0.14 for 13 perceptrons, indicating that the model improves as the number of perceptrons increases (upto a point). As dataset size grows the error rates stabilize and generally decrease, reflecting improved model performance as more data is used for training.

The goal is to train the MLP using maximum likelihood parameter estimation which can be viewed as minimizing the cross entropy loss. So we first identify the optimal number of perceptrons using cross validation. Train the MLP using maximum likelihood (cross entropy loss). Reinitialize the MLP multiple times to avoid getting stuck at local optima & pick the model that achieves the highest log-likelihood on the training data. We then evaluate the model's performance on the test dataset.

After training the MLP with the optimal number of perceptrons for each dataset, the model is evaluated on the test dataset (which contains 100,000 samples). We observe that as the dataset size increases, the error rates on the test dataset tend to decrease. More data helps the model generalize better, reducing the likelihood of overfitting. The model with 100 samples has the highest error rate due to small amount of data for training. The test error rates improves as more data is used for training. For large datasets like 5000 and 10000 samples, the error rates are quite similar (0.17-0.18) indicating that the model has reached a reasonable level of performance and adding more data doesn't significantly change the performance.

Smaller dataset using higher number of perceptrons help the model learn more complex patterns. As datasets grow, fewer perceptrons are preferred to avoid overfitting and maintain a simpler model.

```
1 import numpy as np
2 from sklearn.metrics import accuracy_score
3
4 # Function to apply MAP decision rule and estimate probability of error
5 def map_classification_and_error(mlp, X_test, y_test):
6     # Predict class probabilities using the trained MLP model
7     class_probs = mlp.predict_proba(X_test)
8
9     # For each test sample, select the class with the highest probability (MAP rule)
10    y_pred = np.argmax(class_probs, axis=1)
11
12    # Calculate the probability of error (misclassification rate)
13    error_rate = 1 - accuracy_score(y_test, y_pred)
14
15    return y_pred, error_rate
16
17 # Evaluate each trained MLP on the test set and calculate the error rate
18 for size, (X_train, y_train) in train_datasets.items():
19     best_P = best_perceptrons[size]
20     print(f"Training MLP with {best_P} perceptrons for {size} samples...")
21
22     # Standardize the entire training set and train the model
23     scaler = StandardScaler()
24     X_train_scaled = scaler.fit_transform(X_train)
25
26     # Train the MLP with the optimal number of perceptrons
27     best_mlp = MLPClassifier(hidden_layer_sizes=(best_P,), max_iter=1000, random_state=42, solver='adam', activation='relu')
28     best_mlp.fit(X_train_scaled, y_train)
29
30     # Standardize the test dataset and apply MAP classification and error estimation
31     print("Evaluating MLP on test dataset...")
32     X_test_scaled = scaler.transform(test_X)
33
34     # Apply MAP decision rule and calculate the error rate
35     y_pred, test_error = map_classification_and_error(best_mlp, X_test_scaled, test_y)
36
37     print(f"Test Error Rate for MLP with {best_P} perceptrons: {test_error:.4f}\n")
38
```

Training MLP with 13 perceptrons for 100 samples...  
Evaluating MLP on test dataset...  
Test Error Rate for MLP with 13 perceptrons: 0.1921

Training MLP with 13 perceptrons for 500 samples...  
Evaluating MLP on test dataset...  
Test Error Rate for MLP with 13 perceptrons: 0.1779

Training MLP with 5 perceptrons for 1000 samples...  
Evaluating MLP on test dataset...  
Test Error Rate for MLP with 5 perceptrons: 0.1771

Training MLP with 20 perceptrons for 5000 samples...  
Evaluating MLP on test dataset...  
Test Error Rate for MLP with 20 perceptrons: 0.1718

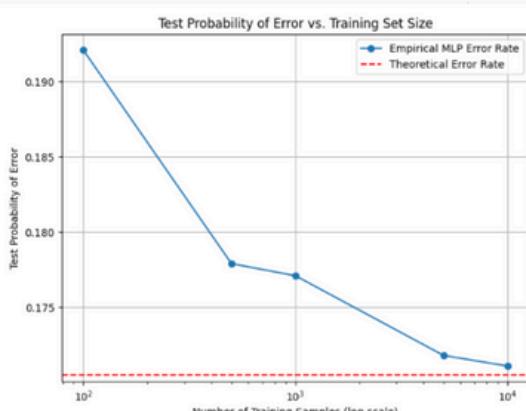
Training MLP with 9 perceptrons for 10000 samples...  
Evaluating MLP on test dataset...  
Test Error Rate for MLP with 9 perceptrons: 0.1711

## 5. Performance Assessment

For each trained MLP, classify the test set and estimate the empirical probability of error using the MAP decision rule. We use the trained MLP classifier as a model to predict class probabilities for each test sample. For each test sample, the MAP decision rule selects the class with the highest posterior probability. After classifying all the test samples, we compare the predicted class labels to the actual class labels in the test set. The error rate is calculated as the fraction of incorrectly classified samples.

Plotting Results: Semilogarithmic plot to show the relationship b/w dataset size & test error rates with the theoretical error as a benchmark

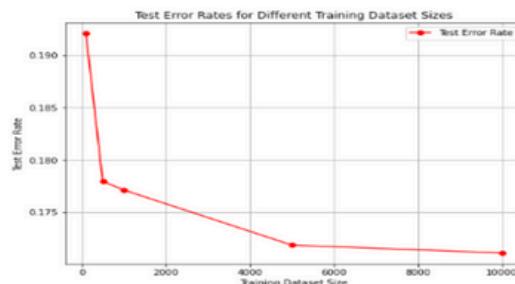
Relationship b/w the training set size and the test probability of error for the MLP model alongside a reference line indicating the theoretical error rate.



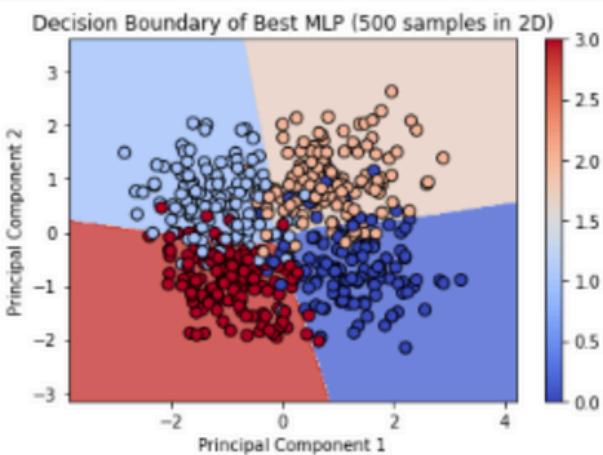
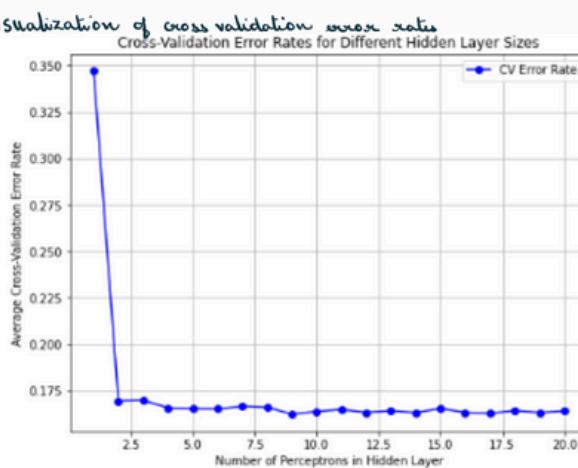
As the number of training samples increases, the test prob. of error decreases suggesting that the model benefits from additional data, likely due to improved generalization as it learns more about the underlying data distribution. The horizontal red-dashed line represents the theoretical error rate which serves as an ideal baseline for model performance. For smaller training sizes the empirical error is above the theoretical error rate indicating that the model struggles to reach the optimal performance. As training set size increases the empirical error rate gets closer to the theoretical limit. This convergence suggests that the model is nearing its best achievable performance. The initial drop in error rate is steep suggesting that the model quickly learns critical patterns. But as we reach larger dataset sizes the rate of improvement slows down.

The graph hints at an asymptotic behavior, where the empirical error rate approaches but doesn't entirely reach the theoretical error rate suggesting that adding more data may yield only marginal gains.

```
Training MLP with 13 perceptrons for 100 samples...
Training MLP on test dataset...
Training MLP with 13 perceptrons for 500 samples...
Training MLP on test dataset...
Training MLP with 5 perceptrons for 1000 samples...
Training MLP on test dataset...
Training MLP with 20 perceptrons for 5000 samples...
Training MLP on test dataset...
Training MLP with 9 perceptrons for 10000 samples...
Training MLP on test dataset...
```



There is a sharp decrease in cross-validation error rates when the no. of perceptrons in the hidden layers increases from 1 to around 3. This suggests that a very small hidden layer captures some critical patterns in the data. After around 3 perceptrons the error rates plateau with only minor fluctuations as the no. of perceptrons increases. This indicates adding more perceptrons beyond this point doesn't provide significant improvement as the model has likely already learned the primary structure of the data.



The plot indicates that the MLP model is reasonably effective at classifying data using complex, non-linear boundaries. However some overlapping suggests that a perfect separation is challenging in this feature space, possibly due to class similarities or projection limitations.

## Software Tools and Implementation

1. NumPy : Used for generating random samples from multivariate normal distributions. The core functionality for dataset generation, random sampling and matrix operations relies on NumPy.
2. Scikit-learn : MLPClassifier is the primary tool used for training the neural network models. It allows us to configure hidden layers, activation functions and optimizers. KFold is used to implement k-fold cross validation to tune the number of perceptions in the hidden layer & estimate error ratio. StandardScaler standardizes the input data to ensure the neural network training is stable and converges more effectively. accuracy\_score is used to compute test error by comparing the predicted labels with the true labels.  
Cross validation ensured that the MLP models were trained on multiple folds of data, mitigating overfitting. By leveraging these tools, we ensured that the training, validation & evaluation steps were conducted systematically and efficiently. The implementation was validated by comparing the MLP's performance against the theoretical optimal classifier, which served as a reference for evaluating the efficacy of the models.

2.

### 1. Setting up the true GMM

A Gaussian Mixture Model is a probabilistic model that assumes data is generated from a mixture of several gaussian distributions. Each component has a mean vector (location of the peak), a covariance matrix (shape and orientation of the gaussian) and a weight (probability of selecting that gaussian component).

The GMM with K components is represented as  $\mu(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$  where K is the number of gaussian components,  $\pi_k$  is the mixing probability of the k<sup>th</sup> component with  $\sum_{k=1}^K \pi_k = 1$ ,  $N(x | \mu_k, \Sigma_k)$  is the gaussian density with mean  $\mu_k$  and covariance  $\Sigma_k$ .

We adjust two of the Gaussian components so they overlap significantly. This can be achieved by setting the distance b/w their mean vectors close to the sum of the average covariance matrix eigenvalues of these two Gaussians. The overlap creates a challenging case for model selection as it may lead to ambiguity in separating the components.

```
1 import numpy as np
2 from sklearn.mixture import GaussianMixture
3
4 # Define the true GMM parameters
5 true_means = np.array([
6     [0, 0],           # Component 1 mean
7     [1, 1],           # Component 2 mean (close to component 1 for overlap)
8     [5, 5],           # Component 3 mean
9     [8, 8]            # Component 4 mean
10 ])
11
12 true_covariances = [
13     [[1, 0.5], [0.5, 1]],   # Component 1 covariance
14     [[1, 0.5], [0.5, 1]],   # Component 2 covariance
15     [[1, -0.3], [-0.3, 1]], # Component 3 covariance
16     [[1, 0.2], [0.2, 1]]    # Component 4 covariance
17 ]
18
19 true_weights = [0.2, 0.3, 0.25, 0.25] # Probabilities of each component
20
```

### 2. Generate Data sets of different sizes

Generate datasets with 10, 100, 1000 samples independently from the true GMM. Each dataset size represents a different data sparsity level. Use the true GMM

to generate independent and identically distributed samples.

```
21 def generate_gmm_data(means, covariances, weights, n_samples):
22     n_components = len(weights)
23     data = []
24     component_choices = np.random.choice(n_components, size=n_samples, p=weights)
25
26     for component in component_choices:
27         point = np.random.multivariate_normal(means[component], covariances[component])
28         data.append(point)
29
30     return np.array(data)
31
32 # Generate datasets
33 dataset_10 = generate_gmm_data(true_means, true_covariances, true_weights, 10)
34 dataset_100 = generate_gmm_data(true_means, true_covariances, true_weights, 100)
35 dataset_1000 = generate_gmm_data(true_means, true_covariances, true_weights, 1000)
36
```

### 3. Cross Validation procedure with different model orders

Evaluate GMMs with 1 to 10 gaussian components. Use the maximum likelihood estimation to fit each candidate GMM. This is achieved with the expectation-maximization algorithm, which iteratively optimizes parameters to maximize the data's log likelihood.

10-fold cross validation splits the data into 10 subsets (folds). For each fold, the model is trained on 9 folds & evaluated on the 10<sup>th</sup> fold. The log likelihoods are averaged over all folds to estimate the model's performance.

```

37 from sklearn.model_selection import KFold
38 from sklearn.mixture import GaussianMixture
39
40 def cross_val_log_likelihood(data, n_components, n_folds=10):
41     kf = KFold(n_splits=n_folds)
42     log_likelihoods = []
43
44     for train_index, test_index in kf.split(data):
45         train_data, test_data = data[train_index], data[test_index]
46         if len(train_data) < n_components:
47             # Skip evaluation if the number of components exceeds available samples in the fold
48             continue
49         gmm = GaussianMixture(n_components=n_components, covariance_type='full')
50         gmm.fit(train_data)
51         log_likelihoods.append(gmm.score(test_data)) # Log-Likelihood on test data
52
53     # Return the mean Log-Likelihood, or a very Low score if we couldn't evaluate any folds
54     return np.mean(log_likelihoods) if log_likelihoods else -np.inf
55
56 # Evaluate model orders from 1 to the minimum of max_components or dataset size
57 def evaluate_model_orders(data, max_components=10, n_folds=10):
58     max_components = min(max_components, len(data)) # Limit by dataset size
59     results = []
60     for n_components in range(1, max_components + 1):
61         avg_log_likelihood = cross_val_log_likelihood(data, n_components, n_folds)
62         results.append(avg_log_likelihood)
63     return results
64
65 # Run evaluation for each dataset
66 results_10 = evaluate_model_orders(dataset_10)
67 results_100 = evaluate_model_orders(dataset_100)
68 results_1000 = evaluate_model_orders(dataset_1000)

```

#### 4. Repeated Experimentation and result aggregation

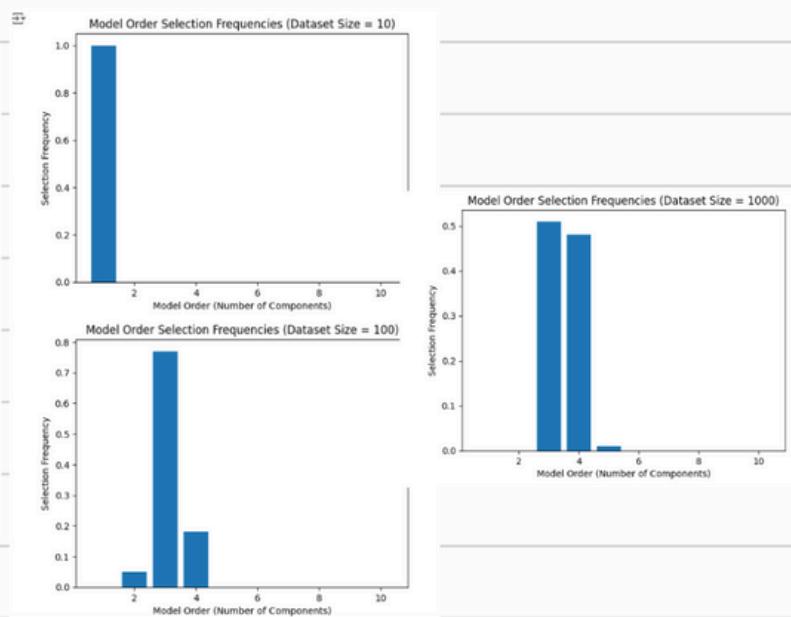
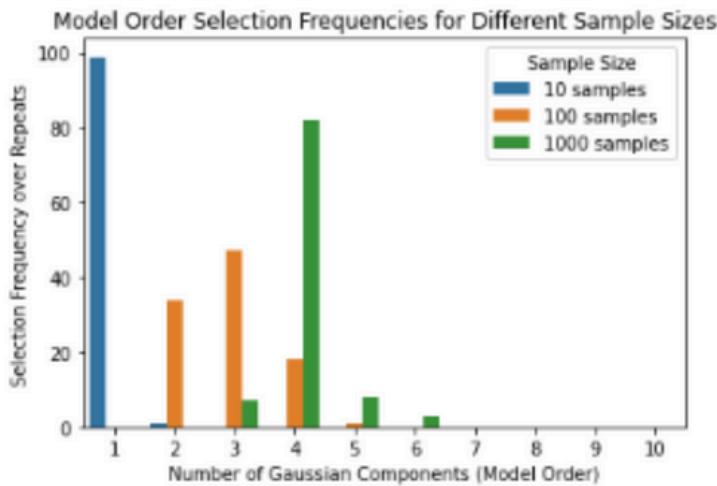
To ensure robust results we repeat the entire cross-validation process 100 times for each dataset size. We record the selected model order with the highest average log-likelihood in each iteration and for each dataset size we determine the frequency of each model order (1 to 10) being selected as the best model across the 100 repetitions.

```

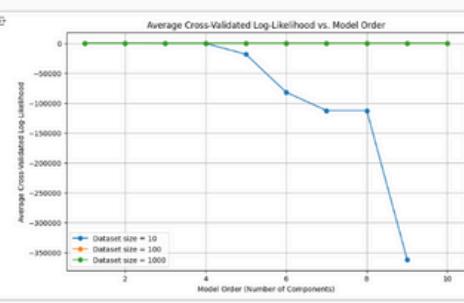
71 def repeat_experiment(data, repetitions=100, max_components=10, n_folds=10):
72     selection_counts = np.zeros(max_components)
73
74     for _ in range(repetitions):
75         log_likelihoods = evaluate_model_orders(data, max_components, n_folds)
76         best_model_order = np.argmax(log_likelihoods) + 1
77         selection_counts[best_model_order - 1] += 1
78
79     return selection_counts / repetitions # Frequency distribution of model order selection
80
81 # Repeat the experiment for each dataset size
82 frequency_10 = repeat_experiment(dataset_10)
83 frequency_100 = repeat_experiment(dataset_100)
84 frequency_1000 = repeat_experiment(dataset_1000)
85
86
87 import matplotlib.pyplot as plt
88
89 def plot_selection_frequencies(frequencies, dataset_size):
90     plt.bar(range(1, len(frequencies) + 1), frequencies)
91     plt.xlabel('Model Order (Number of Components)')
92     plt.ylabel('Selection Frequency')
93     plt.title(f'Model Order Selection Frequencies (Dataset Size = {dataset_size})')
94     plt.show()
95
96 plot_selection_frequencies(frequency_10, 10)
97 plot_selection_frequencies(frequency_100, 100)
98 plot_selection_frequencies(frequency_1000, 1000)
99

```

## RESULTS :



Analysis of the results for model order selection frequencies across different dataset sizes: For small sample size (10 samples) the model consistently selects a 6-MM with only 1 component. This likely indicates that with such a small sample size, the data does not provide enough information to distinguish between multiple underlying gaussian distributions. The model underfits and is unable to capture the complexity of the true 4-component structure. With moderate dataset size of 100, the model can capture more structure in the data and the cross-validation approach gravitates towards selecting 3 components as the best fit followed by some selection of 2 and 4 components. This suggests that with more data the model is able to start identifying a higher number of components but may still not have enough data to consistently recognise all 4 underlying distributions accurately. With 1000 samples the selection frequency peaks sharply at 4 components which matches the true number of components used in the data generation process. The model has sufficient information to correctly identify the true structure of the data.



For the dataset with only 10 samples, the log likelihood steadily decreases as the model order increases. The decline suggests that the model is underfitting when the number of components is increased as it does not have enough data to estimate parameters reliably for higher model orders. The model becomes too complex relative to the limited data, leading to poor generalization.

This indicates that for small datasets, simpler models provide better cross-validated log-likelihood scores as they avoid overfitting. For large datasets, the log likelihood remains high & almost flat across all model orders suggesting that the model consistently achieves good fits. The flat trend indicates that with large datasets the model can reliably capture the underlying structure even as the number of components is varied. The log-likelihood plateau at high values showing that overfitting is less of a concern with large datasets.

## Codes and Results

[https://colab.research.google.com/drive/1rvrflQ0WVLvtmIDyID-UmPUq\\_k7xUya2?usp=sharing](https://colab.research.google.com/drive/1rvrflQ0WVLvtmIDyID-UmPUq_k7xUya2?usp=sharing)