

## Quiz 8

1. 15 points. 1 hour. Using the Wine dataset, use PCA to reduce it to 2 PCA features.

Link: <https://colab.research.google.com/drive/1r8lydpJCCoyEPipXjkAmALvcCK3y5k?usp=sharing>

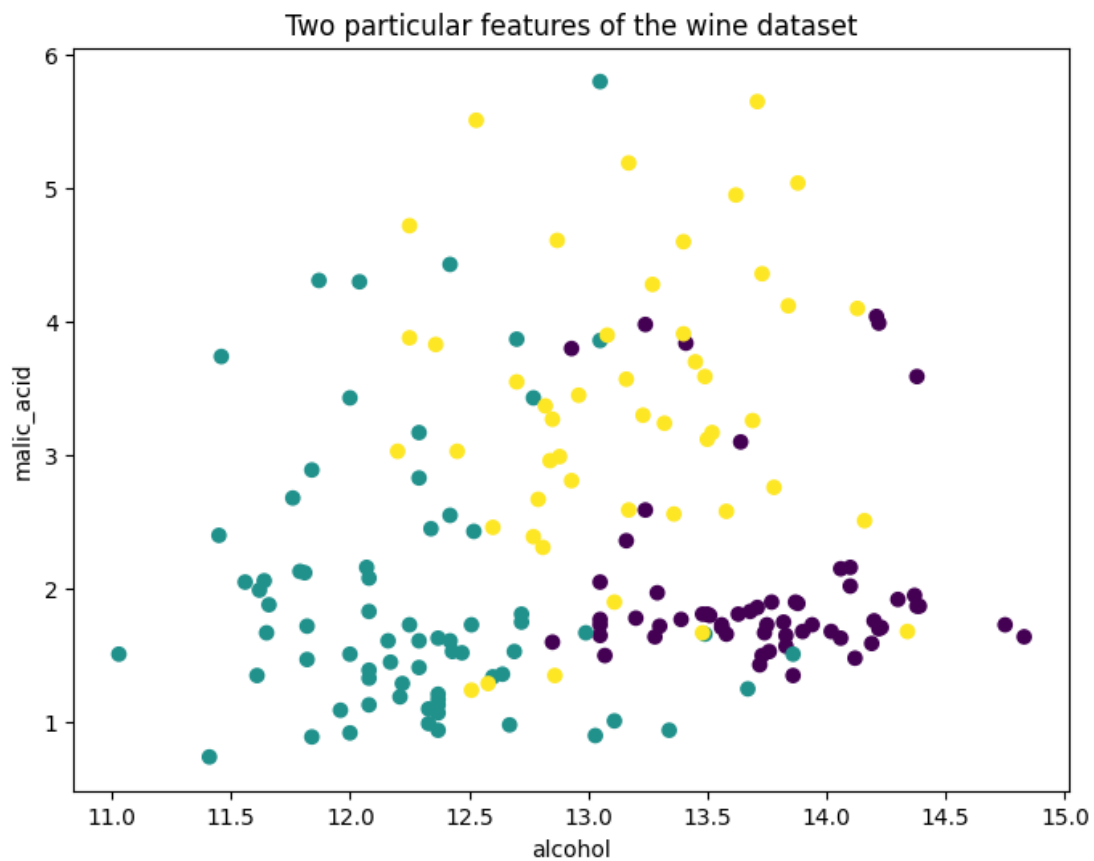
```
[1] 1 #copied mainly from https://machinelearningmastery.com/principal-component-analysis-for-visualization/
2 from sklearn.datasets import load_wine
3 from sklearn.decomposition import PCA
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.pipeline import Pipeline
6 import matplotlib.pyplot as plt
7 import numpy as np
```

```
2 1 # Load dataset
2 2 winedata = load_wine()
3 3 X = winedata['data']
4 4 y = winedata['target']
5 5 print("X shape:", X.shape)
6 6 print("y shape:", y.shape)
7 7 print("Target Output: \n", y) #y has 3 classes
8 8 print("Feature Names: \n", winedata["feature_names"])
9 9 print("Target Names: \n", winedata["target_names"])
10
```

[illegible]

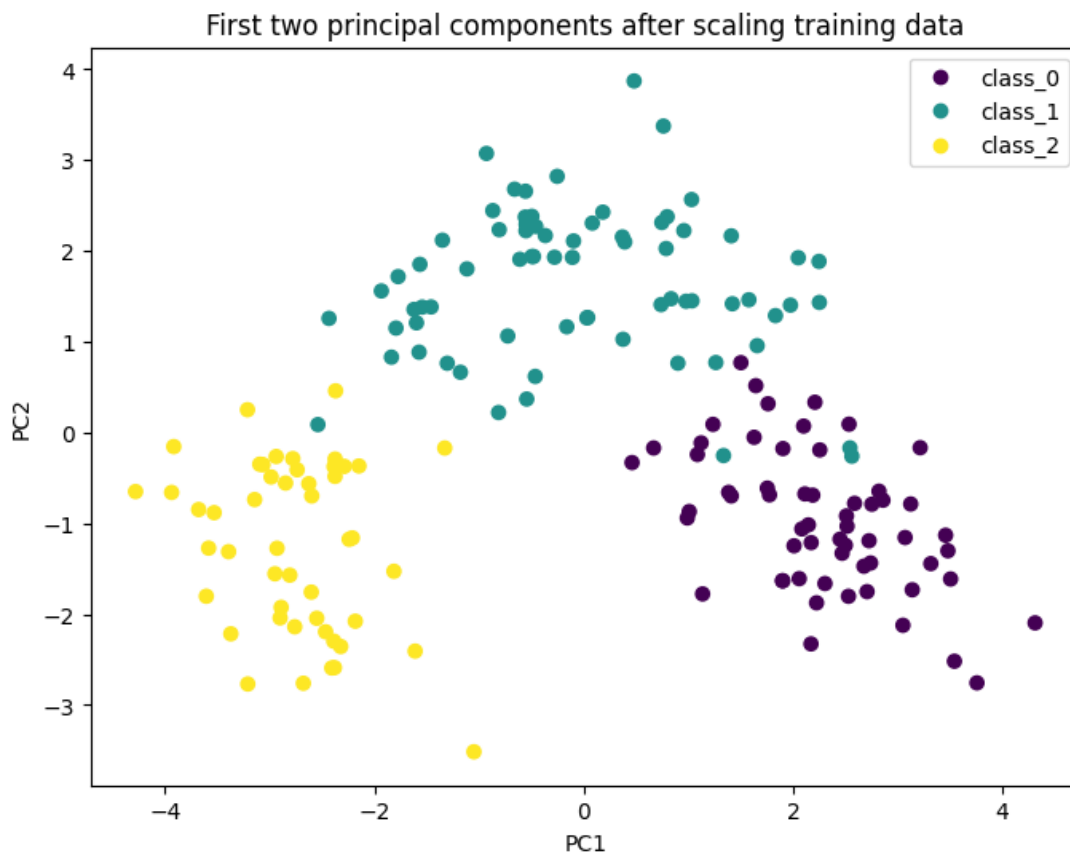
1.1. Visualize the first 2 columns as a scatter plot using all (100%) of the data.

```
[3] 1 # Scatter plot any two features f1 and f2
    2 f1 = 0
    3 f2 = 1
    4 plt.figure(figsize=(8,6))
    5 plt.scatter(X[:,f1], X[:,f2], c=y)
    6 plt.xlabel(winedata["feature_names"][f1])
    7 plt.ylabel(winedata["feature_names"][f2])
    8 plt.title("Two particular features of the wine dataset")
    9 plt.show()
```



1.2. In another graph visualize the scatter plot for the 2 PCA features also using all the data.

```
[4] 1 # Show first two principal components with scaler
2 pca = PCA(n_components=2)
3 pipe = Pipeline([('scaler', StandardScaler()), ('pca', pca)])
4 plt.figure(figsize=(8,6))
5 Xt = pipe.fit_transform(X)
6 plot = plt.scatter(Xt[:,0], Xt[:,1], c=y)
7 plt.legend(handles=plot.legend_elements()[0], labels=list(winedata['target_names']))
8 plt.xlabel("PC1")
9 plt.ylabel("PC2")
10 plt.title("First two principal components after scaling training data")
11 plt.show()
```



1.3. What % of information (“explained variance”) is preserved in the 2 PCA features?

```
[5] 1 # Print the explained variance ratio
2 print("Explained variance ratios:")
3 print(pca.explained_variance_ratio_)
```

```
Explained variance ratios:
[0.36198848 0.1920749 ]
```

```
[6] 1 print("Explained variance of the PCA features: {:.2f}%".format(np.sum(pca.explained_variance_ratio_) * 100))
```

```
Explained variance of the PCA features: 55.41%
```

2. 10 points. 1 hour. Use 25% for testing for the Wine dataset. Plot a graph of the accuracy of the samples when using 1, 2, 3, ..., 13 PCA features along with showing how much % variance is preserved (p%) for each? Use logistic regression.

Link: [https://colab.research.google.com/drive/1PlcOjJH6U9aHv\\_3h9MZE7-ho5EEaXhf?usp=sharing](https://colab.research.google.com/drive/1PlcOjJH6U9aHv_3h9MZE7-ho5EEaXhf?usp=sharing)

```
[1] 1 #copied mainly from https://machinelearningmastery.com/principal-component-analysis-for-visualization/
    2 from sklearn.datasets import load_wine
    3 from sklearn.decomposition import PCA
    4 from sklearn.preprocessing import StandardScaler
    5 from sklearn.pipeline import Pipeline
    6 import matplotlib.pyplot as plt
    7 import numpy as np

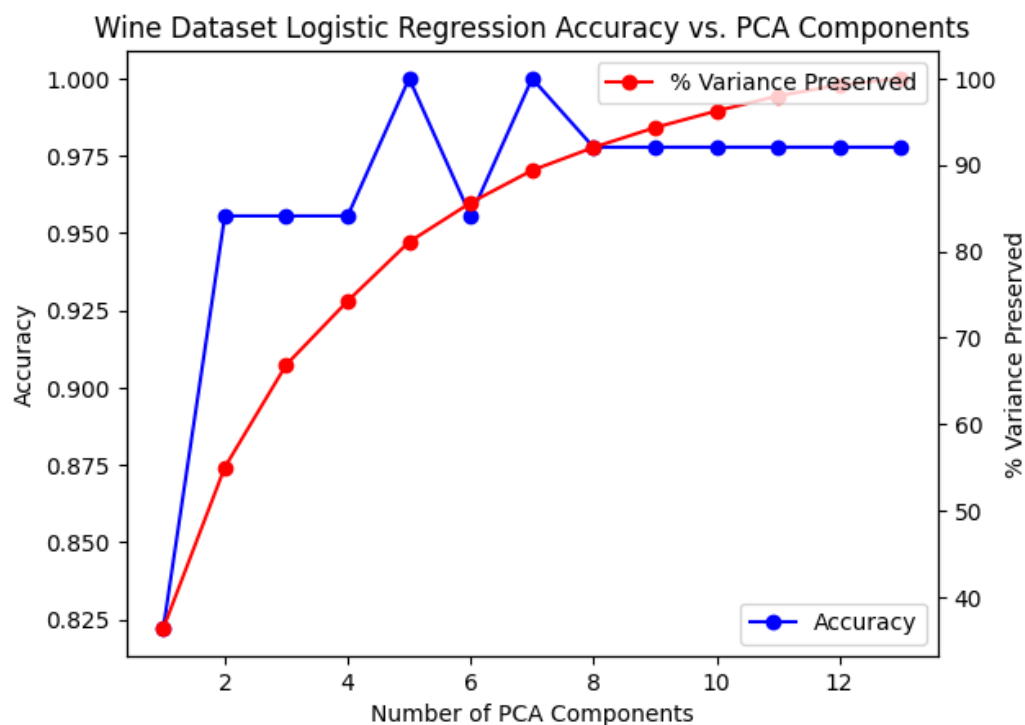
[2] 1 from sklearn.model_selection import train_test_split
    2 from sklearn.linear_model import LogisticRegression

[3] 1 # Load the Wine dataset
    2 wine = load_wine()
    3
    4 # Split the data into training and testing sets
    5 X_train, X_test, y_train, y_test = train_test_split(wine.data, wine.target, test_size=0.25, random_state=42)
    6
    7 # Standardize the features using StandardScaler
    8 scaler = StandardScaler()
    9 X_train_std = scaler.fit_transform(X_train)
   10 X_test_std = scaler.transform(X_test)
   11
   12 # Compute the principal components
   13 pca = PCA()
   14 pca.fit(X_train_std)
   15
   16 # Compute the variance explained by each component
   17 var_exp = pca.explained_variance_ratio_
   18
   19 # Compute the cumulative variance explained
   20 cum_var_exp = np.cumsum(var_exp)
   21
   22 # Initialize lists to store the accuracies and % variance preserved
   23 accuracies = []
   24 percent_var_preserved = []
```

```

[4] 1 # Loop over the number of PCA components
2 for i in range(1, 14):
3     # Compute the first i principal components
4     pca = PCA(n_components=i)
5     pca.fit(X_train_std)
6     X_train_pca = pca.transform(X_train_std)
7     X_test_pca = pca.transform(X_test_std)
8
9     # Fit a logistic regression model to the transformed data
10    lr = LogisticRegression(random_state=42)
11    lr.fit(X_train_pca, y_train)
12
13    # Compute the accuracy on the transformed testing data
14    accuracy = lr.score(X_test_pca, y_test)
15    accuracies.append(accuracy)
16
17    # Compute the % variance preserved
18    percent_var = cum_var_exp[i-1]*100
19    percent_var_preserved.append(percent_var)
20
21 # Plot the accuracy vs. the number of PCA components
22 plt.plot(range(1, 14), accuracies, 'bo-', label='Accuracy')
23 plt.xlabel('Number of PCA Components')
24 plt.ylabel('Accuracy')
25 plt.title('Wine Dataset Logistic Regression Accuracy vs. PCA Components')
26 plt.legend(loc='lower right')
27
28 # Plot the % variance preserved vs. the number of PCA components
29 plt.twinx()
30 plt.plot(range(1, 14), percent_var_preserved, 'ro-', label='% Variance Preserved')
31 plt.ylabel('% Variance Preserved')
32 plt.legend(loc='upper right')
33
34 plt.show()

```



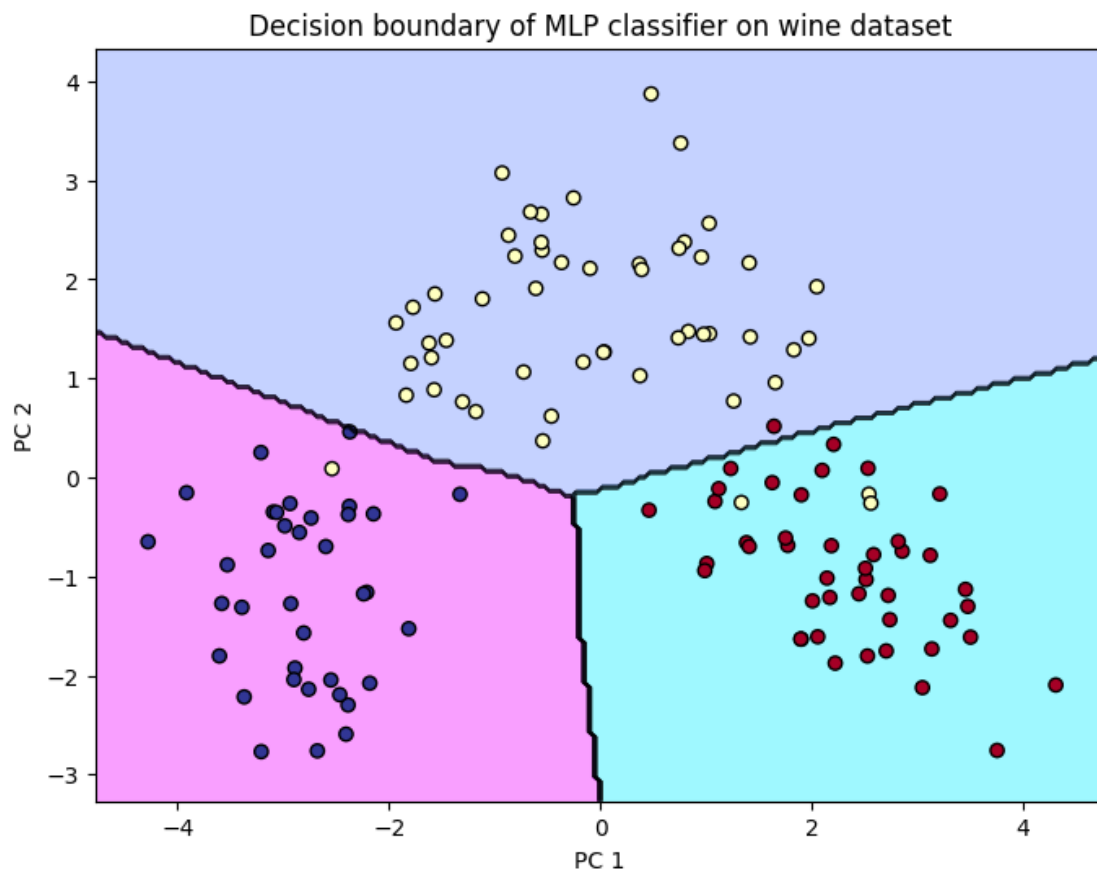
3. 15 points. 2 hours. Decision boundary. Plot the decision boundary for PCA-reduced 2 dimensions of the wine dataset. Use ML Perceptron in <https://python-course.eu/machine-learning/neural-networks-with-scikit.php> with 2 hidden layers.

Link: <https://colab.research.google.com/drive/1tmF-MHHgtq0zERJgn0RysYCx30RaVEEN?usp=sharing>

```
[1] 1 # Load necessary libraries and modules
    2 import numpy as np
    3 import matplotlib.pyplot as plt
    4 from sklearn.datasets import load_wine
    5 from sklearn.model_selection import train_test_split
    6 from sklearn.neural_network import MLPClassifier
    7 from sklearn.decomposition import PCA
    8 from sklearn.preprocessing import StandardScaler
    9 from sklearn.pipeline import Pipeline

[2] 1 # Load the wine dataset
    2 wine_data = load_wine()
    3 X = wine_data['data']
    4 y = wine_data['target']
    5
    6 # Perform PCA on the training data
    7 pca = PCA(n_components=2)
    8 scaler = StandardScaler()
    9 X_pca = pca.fit_transform(scaler.fit_transform(X))
   10
   11 # Split the data into training and testing sets
   12 X_train, X_test, y_train, y_test = train_test_split(X_pca, y, test_size=0.25, random_state=42)
   13
   14 # Train a MLP classifier with 2 hidden layers
   15 model = MLPClassifier(hidden_layer_sizes=(10, 10), max_iter=1000)
   16 model.fit(X_train, y_train)
   17
   18 # Plot the decision boundary
   19 h = 0.05
   20 x_min, x_max = X_train[:, 0].min() - 0.5, X_train[:, 0].max() + 0.5
   21 y_min, y_max = X_train[:, 1].min() - 0.5, X_train[:, 1].max() + 0.5
   22 xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
   23 Z = model.predict(np.c_[xx.ravel(), yy.ravel()])
   24 Z = Z.reshape(xx.shape)
```

```
[3] 1 # Create the plot
2 plt.figure(figsize=(8, 6))
3 plt.contourf(xx, yy, Z, cmap=plt.cm.cool, alpha=0.4)
4 plt.contour(xx, yy, Z, colors='black', linewidths=0.5)
5 plt.scatter(X_train[:, 0], X_train[:, 1], c=y_train, cmap=plt.cm.RdYlBu, edgecolors='black')
6 plt.xlabel('PC 1')
7 plt.ylabel('PC 2')
8 plt.title('Decision boundary of MLP classifier on wine dataset')
9 plt.show()
```



4. Using the MNIST data of 784 dimensions (28 x 28 pixels).

Link: <https://colab.research.google.com/drive/1wVTUEd6ipaYL4iOpcUoVsBDmXxUPOkHJ?usp=sharing>

- 4.1. 5 points. 0.5 hours. Show Quiz 7 problem #3 here again as a benchmark. Using 50,000 data samples for training and 20,000 for testing, apply standard scalar transform from the training dataset to both your training and testing datasets. Report the accuracy and print out the confusion matrix.

```
[25] 1 import numpy as np
      2 import matplotlib.pyplot as plt
      3 from sklearn.datasets import fetch_openml
      4 from sklearn.neural_network import MLPClassifier
      5 from sklearn.model_selection import train_test_split
      6 from sklearn.metrics import accuracy_score, confusion_matrix
      7 from sklearn import metrics
      8 from sklearn.preprocessing import StandardScaler
```

```
[26] 1 # Load the MNIST dataset
      2 mnist = fetch_openml("mnist_784")
      3
      4 # Split the data into training and testing sets
      5 X_train, X_test, y_train, y_test = train_test_split(mnist.data, mnist.target, train_size=50000, test_size=20000)
      6
      7 # Scale the data to be between 0 and 1
      8 X_train = X_train / 255.0
      9 X_test = X_test / 255.0
```

/usr/local/lib/python3.9/dist-packages/sklearn/datasets/\_openml.py:968: FutureWarning: The default value of `parser` will c  
warn(

```
[27] 1 # Initialize the MLPClassifier
      2 clf = MLPClassifier(hidden_layer_sizes=(100, 100), max_iter=10)
      3
      4 # Train the MLPClassifier on the training data
      5 clf.fit(X_train, y_train)
      6
      7 # Use the trained MLPClassifier to make predictions on the testing data
      8 y_pred = clf.predict(X_test)
      9
     10 # Scale the training set
     11 scaler = StandardScaler()
     12 X_train_scaled = scaler.fit_transform(X_train)
     13
     14 # Scale the testing set using the same scaler
     15 X_test_scaled = scaler.transform(X_test)
```



```

[28] 1 # Initialize the MLPClassifier
2 clf = MLPClassifier(hidden_layer_sizes=(100, 100), max_iter=10)
3
4 # Train the MLPClassifier on the training data
5 clf.fit(X_train_scaled, y_train)
6
7 # Use the trained MLPClassifier to make predictions on the testing data
8 y_pred = clf.predict(X_test_scaled)
9
10 # Print the accuracy score of the MLPClassifier on the testing data
11 accuracy = accuracy_score(y_test, y_pred)
12 print('Accuracy:', accuracy)
13
14 # Show the confusion matrix of the MLPClassifier on the testing data
15 print("\nConfusion Matrix:")
16 confusion_matrix(y_test, y_pred)

```

/usr/local/lib/python3.9/dist-packages/sklearn/neural\_network/\_multilayer\_perceptron.py:686: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (10) reached without convergence. The results might be not optimal.  
 warnings.warn(  
 Accuracy: 0.9679

Confusion Matrix:

```

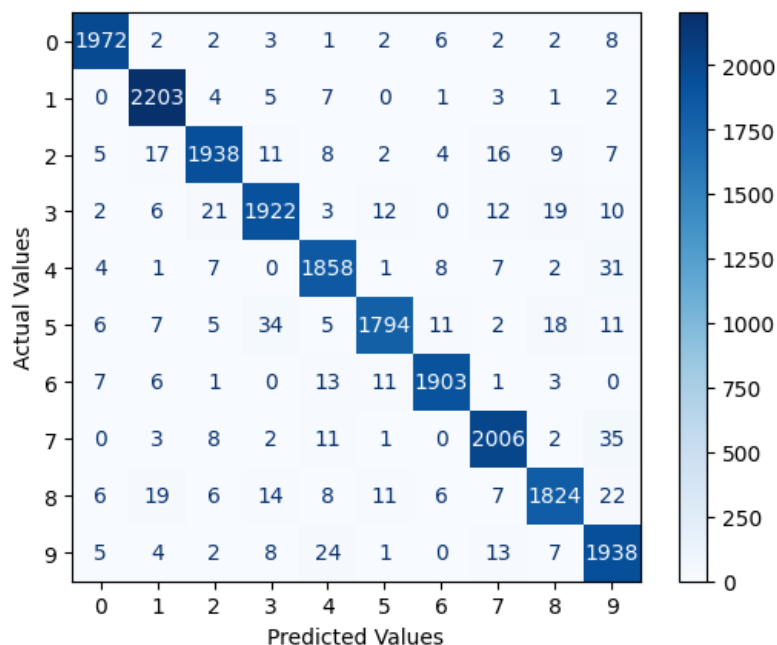
array([[1972,  2,  2,  3,  1,  2,  6,  2,  2,  8],
       [ 0, 2203,  4,  5,  7,  0,  1,  3,  1,  2],
       [ 5, 17, 1938, 11,  8,  2,  4, 16,  9,  7],
       [ 2,  6, 21, 1922,  3, 12,  0, 12, 19, 10],
       [ 4,  1,  7,  0, 1858,  1,  8,  7,  2, 31],
       [ 6,  7,  5, 34,  5, 1794, 11,  2, 18, 11],
       [ 7,  6,  1,  0, 13, 11, 1903,  1,  3,  0],
       [ 0,  3,  8,  2, 11,  1,  0, 2006,  2, 35],
       [ 6, 19,  6, 14,  8, 11,  6,  7, 1824, 22],
       [ 5,  4,  2,  8, 24,  1,  0, 13,  7, 1938]])

```

```

[29] 1 # Confusion matrix
2 cm = metrics.confusion_matrix(y_test, y_pred)
3 classes = np.unique(mnist.target)
4 cm_display = metrics.ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=classes)
5 cm_display.plot(cmap=plt.cm.Blues)
6 plt.xlabel('Predicted Values')
7 plt.ylabel('Actual Values')
8 plt.show()

```



- 4.2. 20 points. 2 hours. Using the same neural network model as in 4.1 above, but using the PCA features that preserve up to 95% of the information content, find a new training and testing dataset using PCA transformation based on the training dataset. Apply standard scalar transform from the training dataset to both your training and testing datasets. Report the accuracy and print out the confusion matrix for PCA with reduced features.

```
[41] 1 import numpy as np
      2 import matplotlib.pyplot as plt
      3 from sklearn.datasets import fetch_openml
      4 from sklearn.neural_network import MLPClassifier
      5 from sklearn.model_selection import train_test_split
      6 from sklearn.metrics import accuracy_score, confusion_matrix
      7 from sklearn import metrics
      8 from sklearn.preprocessing import StandardScaler
      9 from sklearn.decomposition import PCA
     10 from sklearn.pipeline import Pipeline
```

```
[42] 1 # Load the MNIST dataset
      2 X, y = fetch_openml("mnist_784", version=1, return_X_y=True)
      3
      4 # Normalize intensity of images to make it in the range [0,1] since 255 is the max (white).
      5 X = X / 255.0
      6
      7 # Split the data into train/test sets
      8 X_train, X_test = X[:50000], X[50000:]
      9 y_train, y_test = y[:50000], y[50000:]
```

/usr/local/lib/python3.9/dist-packages/sklearn/datasets/\_openml.py:968: FutureWarning: The default value of `parser` will change from warn(

```
[43] 1 # Apply PCA to training set with 95% variance preserved
      2 pca = PCA(n_components=0.95)
      3 pipe = Pipeline([('scaler', StandardScaler()), ('pca', pca)])
      4 X_train_pca = pipe.fit_transform(X_train)
      5 X_test_pca = pipe.transform(X_test)
      6 X_train_pca = pca.fit_transform(X_train)
      7
      8 # Apply PCA to testing set with the same transformation
      9 X_test_pca = pca.transform(X_test)
     10
     11 # Train neural network model
     12 classifier = MLPClassifier(hidden_layer_sizes=(50,), max_iter=10, alpha=1e-4,
     13 | | | solver='sgd', verbose=10, random_state=1,
     14 | | | learning_rate_init=.1)
     15
     16 clf = classifier.fit(X_train_pca, y_train)
     17
     18 # Make predictions on test data
     19 y_pred = clf.predict(X_test_pca)
```

```

↳ Iteration 1, loss = 0.32053499
Iteration 2, loss = 0.14439844
Iteration 3, loss = 0.11033504
Iteration 4, loss = 0.09065064
Iteration 5, loss = 0.07793818
Iteration 6, loss = 0.06816407
Iteration 7, loss = 0.05963315
Iteration 8, loss = 0.05319021
Iteration 9, loss = 0.04850037
Iteration 10, loss = 0.04360186
/usr/local/lib/python3.9/dist-packages/sklearn/neural_network/_multilayer_perceptron.py:686: ConvergenceWarning: Stochastic Optimizer:
warnings.warn(

```

```

[44] 1 # Print the accuracy score of the MLPClassifier on the testing data
2 accuracy = accuracy_score(y_test, y_pred)
3 print('Accuracy:', accuracy)
4
5 # Show the confusion matrix of the MLPClassifier on the testing data
6 print("\nConfusion Matrix:")
7 confusion_matrix(y_test, y_pred)

```

Accuracy: 0.9833

Confusion Matrix:

```

array([[4879, 1, 10, 2, 2, 1, 4, 3, 2, 5],
       [ 0, 5565, 6, 4, 2, 1, 7, 8, 3, 0],
       [11, 13, 4969, 8, 13, 0, 4, 27, 14, 2],
       [ 4, 2, 25, 4983, 0, 16, 0, 7, 15, 13],
       [ 4, 13, 2, 0, 4789, 0, 11, 6, 2, 52],
       [13, 5, 3, 51, 5, 4411, 24, 3, 9, 14],
       [15, 4, 3, 2, 10, 6, 4861, 1, 3, 0],
       [ 4, 19, 15, 8, 7, 3, 0, 5118, 2, 24],
       [15, 16, 16, 42, 2, 13, 14, 7, 4760, 18],
       [ 9, 8, 2, 21, 24, 16, 0, 30, 4, 4830]])

```

```

[45] 1 # Confusion matrix
2 cm = metrics.confusion_matrix(y_test, y_pred)
3 classes = np.unique(mnist.target)
4 cm_display = metrics.ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=classes)
5 cm_display.plot(cmap=plt.cm.Blues)
6 plt.xlabel('Predicted Values')
7 plt.ylabel('Actual Values')
8 plt.show()

```

