

# MACHINE LEARNING (COMP7703)

Assignment - Semester 1, 2025.

**Student Name:** Volter Entoma

**Student ID:** 44782711

## Machine Learning Approaches for Height Classification from Table Tennis Swing Data: A Comparative Study

### Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Aims</b>	<b>3</b>
<b>3</b>	<b>Data</b>	<b>3</b>
3.1	Data split . . . . .	4
3.2	Data Imbalance . . . . .	4
3.3	Data Standardisation . . . . .	4
<b>4</b>	<b>Principle Component Analysis</b>	<b>5</b>
<b>5</b>	<b>K-NN Classifier</b>	<b>7</b>
5.1	Accuracy of K-NN Classifier . . . . .	8
5.2	K-Fold Cross-validation . . . . .	9
5.3	ROC Curve . . . . .	10
5.4	Confusion Matrix . . . . .	11
5.5	Precision, Recall, and F1 score . . . . .	12
<b>6</b>	<b>Random Forest Classifier</b>	<b>13</b>
6.1	Accuracy of Random Forest Classifier . . . . .	13
6.1.1	Computational time after dimensionality reduction . . . . .	14
6.2	K-Fold Cross-validation . . . . .	15
6.3	ROC curve . . . . .	16
6.4	Confusion matrix . . . . .	17
6.5	Precision, Recall, and F1 score . . . . .	17
<b>7</b>	<b>Deep Neural Networks</b>	<b>18</b>
7.1	Training . . . . .	18
7.1.1	Effect of hidden layer sizes . . . . .	19
7.1.2	Training 192-192 DNN model . . . . .	19
7.2	ROC Curve . . . . .	20
7.3	Confusion Matrix . . . . .	21
7.4	Precision, Recall, and F1 score . . . . .	21

<b>8</b>	<b>Conclusion</b>	<b>22</b>
<b>9</b>	<b>Appendix</b>	<b>23</b>
9.1	All Project Code . . . . .	23

# 1 Introduction

There are many different machine learning algorithms, and many metrics to evaluate them. In this assignment, we will focus on the K-Nearest Neighbour (K-NN) algorithm, Random Forest, and Deep Neural Networks, and the accuracy of each model. We will also look at the different metrics to evaluate the performance of each model, and use these metrics to compare the suitability of each model for the given data.

## 2 Aims

The primary aim of this paper is to evaluate and compare the performance of three machine learning classifiers—K-Nearest Neighbours (K-NN), Random Forest, and Deep Neural Networks (DNN)—on the TTSWING dataset for the task of classifying player height categories based on swing sensor data. The specific objectives are as follows:

- To preprocess and analyze the TTSWING dataset, addressing issues such as data imbalance and feature scaling.
- To implement and optimize K-NN, Random Forest, and DNN classifiers for the height classification task.
- To investigate the impact of dimensionality reduction (via PCA) on model performance and computational efficiency.
- To assess each model using a range of evaluation metrics, including accuracy, ROC curve, confusion matrix, precision, recall, and F1 score.
- To perform cross-validation to evaluate model generalizability and prevent overfitting.
- To compare the strengths, weaknesses, and trade-offs of each classifier in terms of predictive performance, and computational cost
- To propose potential improvements and future directions for each classification strategy.

## 3 Data

The data used in this assignment is the TTSWING dataset, which is the provided dataset for this assignment. The dataset contains the 34 features, which are the 34 different measurements of table tennis swings. The dataset also contains the height of the player, which is the target variable. The dataset is a CSV file, and can be read into Python using the pandas library. The dataset is split into three categories: high, medium, and low.

### 3.1 Data split

The data is split into three height categories: high, medium, and low. The data is split into 60% training data, 20% testing data, and 20% validation data. The training data is used to train the each model, the testing data is used to test the accuracy of each model, and the validation data is used to compare the performance of the models.

### 3.2 Data Imbalance

After splitting the data, there is an imbalance in the amount of training data assigned to each category in height. There are 17364 (29.73%) data points in the high, 40800 (42.03%) data points in the medium, and 27450 (28.23%) in the low category. This is a significant imbalance; this imbalance is problematic because many machine learning algorithms tend to be biased toward the majority class. This can result in poor predictive performance for the minority classes, as the model may learn to ignore them in favor of achieving higher overall accuracy. In imbalanced datasets, accuracy becomes a misleading metric, since a model can achieve high accuracy by simply predicting the majority class most of the time, while failing to correctly classify minority class instances. This is especially concerning when the minority classes are of particular interest or importance.

The data imbalanced is addressed by the use of Synthetic Minority Over-sampling Technique (SMOTE) to generate synthetic data points for the low category. The SMOTE algorithm generates synthetic data points by interpolating between existing data points in the low category. This is done by selecting a random data point from the low category, and then selecting a random data point from the  $k$  nearest neighbours of that data point. The synthetic data point is then generated by taking a weighted average of the two data points.

After applying SMOTE, the training data is balanced, with 24552 data points in each category. All subsequent analysis is performed on the balanced training data. The testing and validation data is not balanced, and is used to test the accuracy of the models.

### 3.3 Data Standardisation

The data is standardized using the StandardScaler from the `sklearn` library. The StandardScaler standardizes the data by removing the mean and scaling to unit variance. This is done to ensure that all features have the same scale, and to improve the performance of the models.

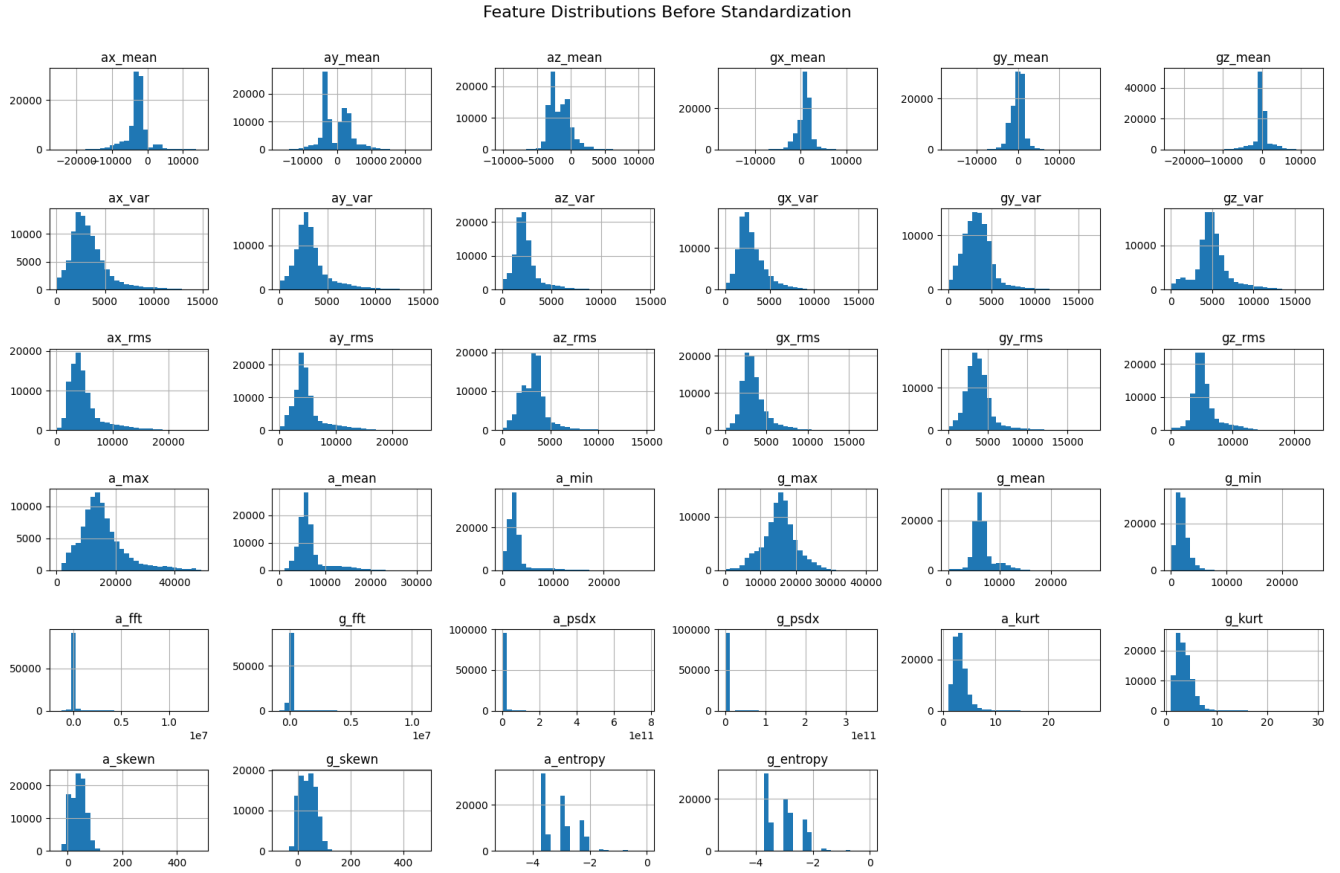


Figure 1: Distribution of the data before standardization.

## 4 Principle Component Analysis

To investigate the effect of dimensionality reduction on the accuracy of the models, we will use Principle Component Analysis (PCA) to reduce the dimensionality of the data. PCA is a linear transformation that transforms the data into a new coordinate system, where the first coordinate is the direction of maximum variance, the second coordinate is the direction of maximum variance orthogonal to the first coordinate, and so on.

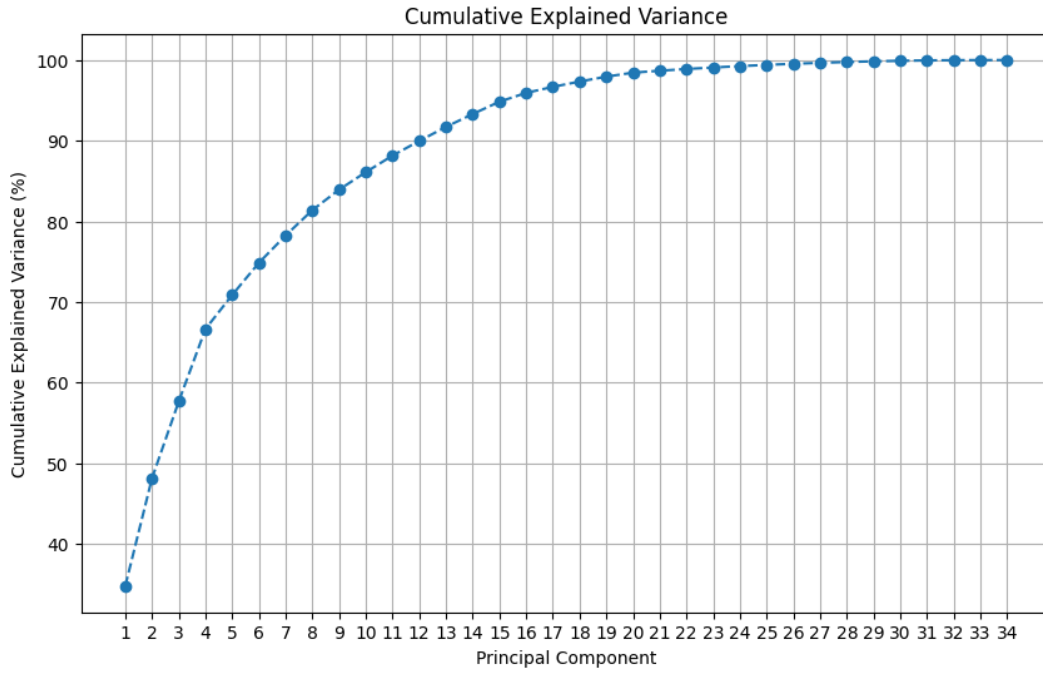


Figure 2: Cumulative variance explained by each principal component.

As shown in Figure 2, approximately 95% of the total variance is explained by the top 16 principal components. This allows us to reduce the dimensionality of our data from 34 dimensions to 16 dimensions while retaining the vast majority of the information content.

By taking a look at the correlation matrix of the data, as shown in Figure 3, we can see that there are many features that are highly correlated with each other. These off-diagonal correlations could explain why the PCA transformation is able to reduce the dimensionality of the data while retaining a high percentage of the variance.

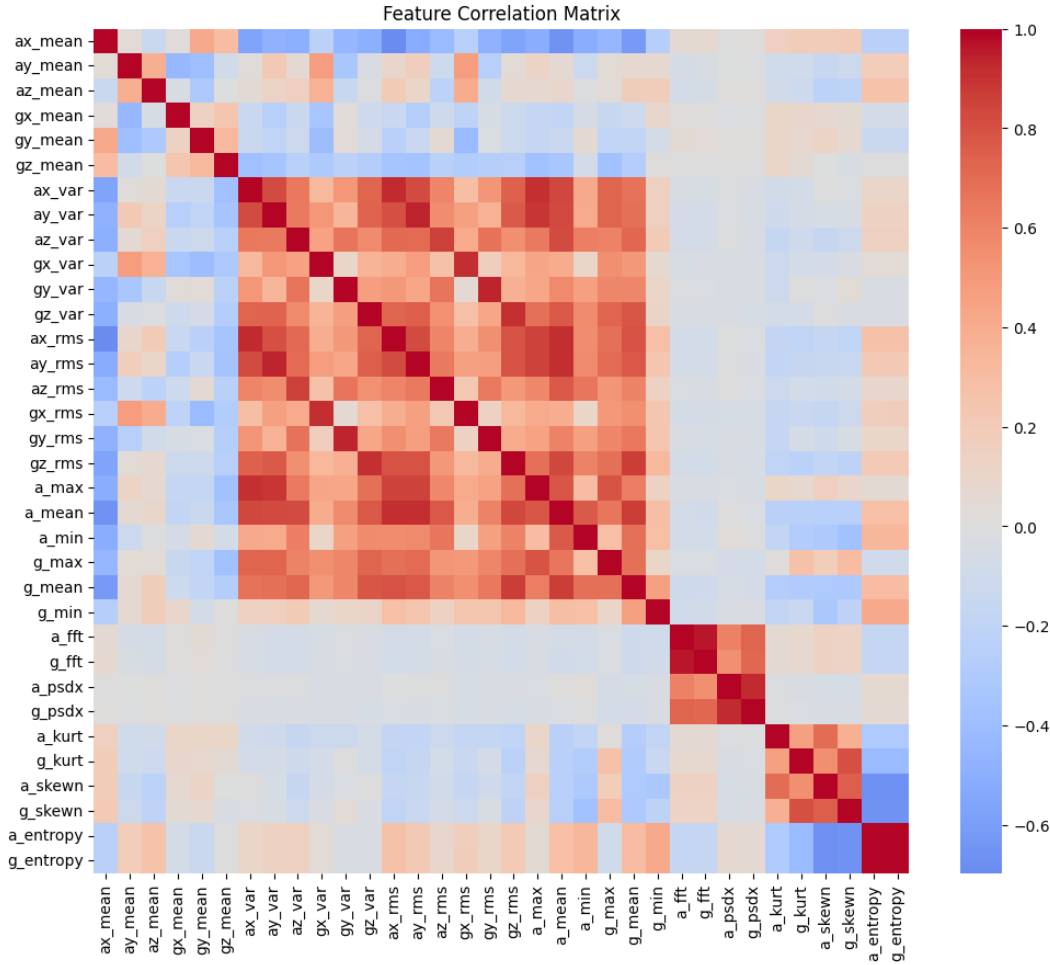


Figure 3: Correlation matrix of the standardized data.

Dimensional reduction offers several benefits:

- Reduced computational complexity in subsequent modeling
- Mitigation of the curse of dimensionality
- Removal of potentially noisy or redundant features
- Improved model interpretability

The effectiveness of this dimensionality reduction will be evaluated by comparing model performance on both the original and PCA-reduced datasets.

## 5 K-NN Classifier

The K-Nearest Neighbour (K-NN) algorithm is a simple and effective machine learning algorithm that can be used for both classification and regression tasks. The K-NN algorithm works by finding the  $k$  nearest neighbours of a data point, and then predicting the class of the data point based on the classes

of the  $k$  nearest neighbours. The K-NN algorithm is a non-parametric algorithm, which means that it does not make any assumptions about the distribution of the data. This makes the K-NN algorithm very flexible, and it can be used for a wide variety of tasks, including for our case of the classifying the height of a table tennis player based on their swing data.

The K-NN algorithm was applied to the data to classify the height of the player based on their swing data.

## 5.1 Accuracy of K-NN Classifier

Using the data split outlined in Section 3.1, the K-NN classifier was trained on the training data, and then tested on the testing data.

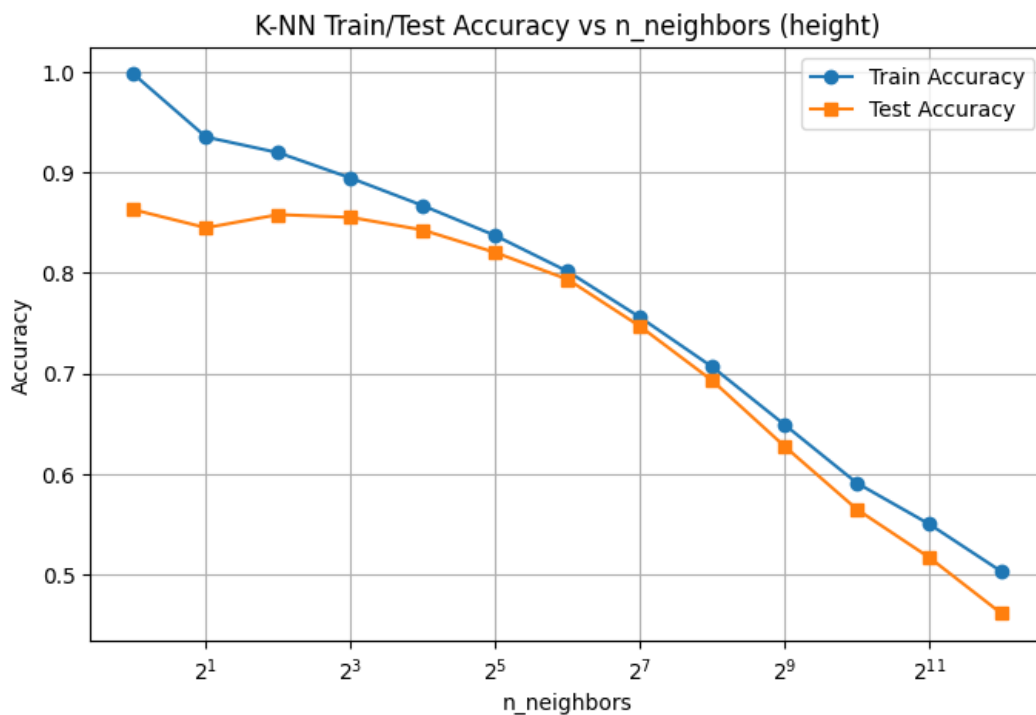


Figure 4: Accuracy of K-NN classifier against the train and test data.

The accuracy of the K-NN classifier is heavily dependent on the number of neighbours ( $k$ ) used in the algorithm. Figure 4 shows the accuracy of the K-NN classifier reaches a maximum at 1 neighbour, and then decreases as the number of neighbours increases. This could be related to the high dimensionality of the data.

The K-NN classifier is particularly sensitive to the curse of dimensionality; as the number of dimensions increases, point within the hyperspace are more likely to be equidistant from each other. This means that the K-NN algorithm is less able to distinguish between points, and hence the accuracy of the K-NN classifier decreases. This is a common problem with high dimensional data, and is one of the reasons why dimensionality reduction techniques such as PCA are used.

We can confirm this by looking at the accuracy of the K-NN classifier on the PCA-reduced data.



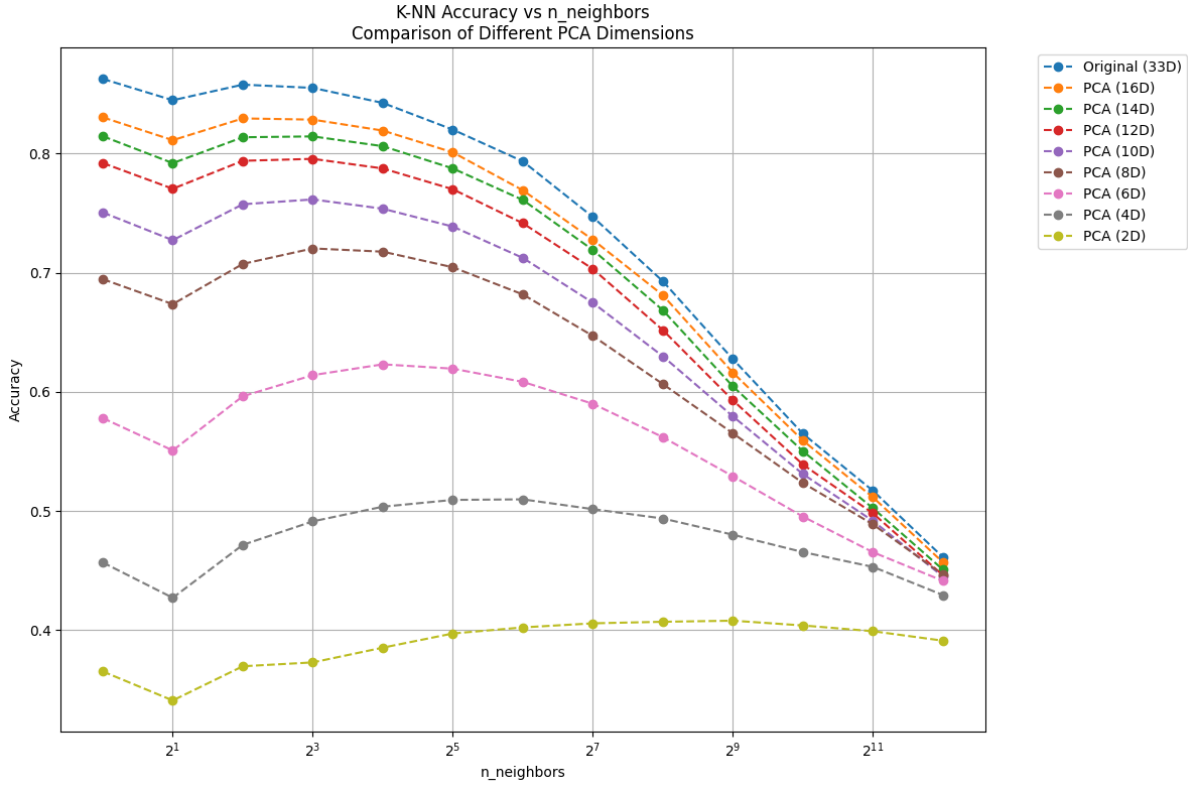


Figure 5: Accuracy of K-NN classifier on PCA-reduced data against the test data.

The accuracy of the K-NN classifier on the PCA-reduced data is shown in Figure 5. The curve for reduced data to 2-dimensions shows the best accuracy occurs with  $2^9$  neighbour, which suggests that the points become more "spread-out" in 2-dimensions and the model is able to distinguish between different categories better as the number of neighbours increases.

However, despite the dimensionality reduction, the accuracy of the K-NN classifier is still highest at 1 neighbour, and then decreases as the number of neighbours increases. This suggests that the decrease in accuracy as neighbours increase is not solely due to the high dimensionality of the data, but also due to the nature of the K-NN algorithm on this particular dataset. This idea is supported by the fact that the curve for the PCA-reduced data (to 16 dimensions) is similar to the curve for the original data.

Additionally, the accuracy of the K-NN classifier decreases as the number of dimensions decreases. If the curse of dimensionality was a significant problem for this dataset, we would expect the accuracy to increase. This suggests that the curse of dimensionality is not a significant problem for this dataset, and that decreasing the number dimensions is not beneficial to the model, as it loses too much information.

## 5.2 K-Fold Cross-validation

The best accuracy of the K-NN classifier is achieved with one neighbour, and the accuracy is approximately 86.26%. However, one neighbour could be overfitting the data, therefore is important to explore different sets of training and testing data, to see if the accuracy is consistent across different sets of data. This is done by using K-Fold Cross-validation, which splits the data into K different sets, and then trains and tests the model on each set. The accuracy of the model is then averaged over all K sets.

We use 10-fold cross-validation, which means that the data is split into 10 different sets, and then the

model is trained and tested on each set. Note that the choice of the number of folds was arbitrary, and was chosen to be 10 because it is a common choice in the literature.

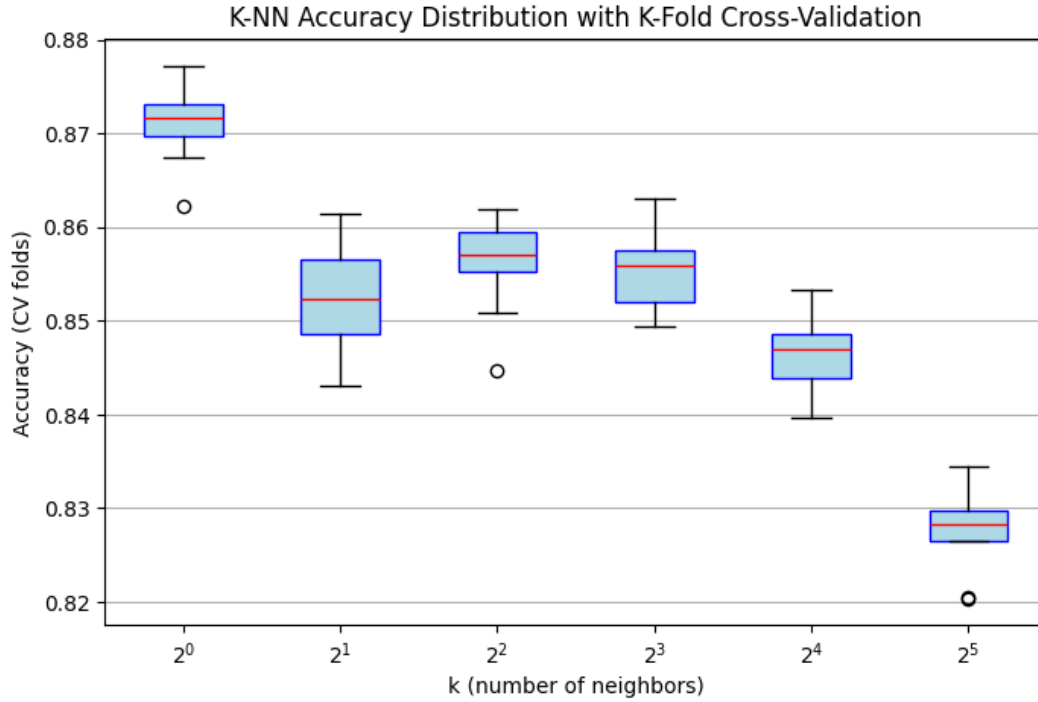


Figure 6: Accuracy of K-NN classifier with K-Fold Cross-validation. Each box and whisker plots shows the spread of the accuracy against the test data across the 10 folds. The red line shows the mean accuracy across the 10 folds.

The accuracy of the K-NN classifier with K-Fold cross-validation is shown in Figure 6. With 10-fold cross-validation, the highest average accuracy occurs with one neighbour. Notably, the minimum accuracy across the 10 folds for one neighbour is approximately equal to the maximum accuracy observed for two, four, and eight neighbours. This suggests that one neighbour is indeed the optimal choice for this dataset.

Additionally, the variance in accuracy across folds for one neighbour is comparable to the variance observed for other values of  $k$ . If one neighbour were overfitting the data, we would expect to see significantly higher variance across the folds for  $k=1$  compared to other values of  $k$ . Since this is not the case, we can conclude that, for this dataset, the K-NN classifier with one neighbour does not appear to overfit, as evidenced by the consistently low variance.

### 5.3 ROC Curve

We use the  $k=1$  K-NN model with the highest accuracy in the 10-fold cross-validation to plot the ROC curve. The ROC curve is a graphical representation of the performance of a binary classifier as the discrimination threshold is varied. It plots the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings. The area under the ROC curve (AUC) is a single scalar value that summarizes the performance of the classifier across all thresholds.

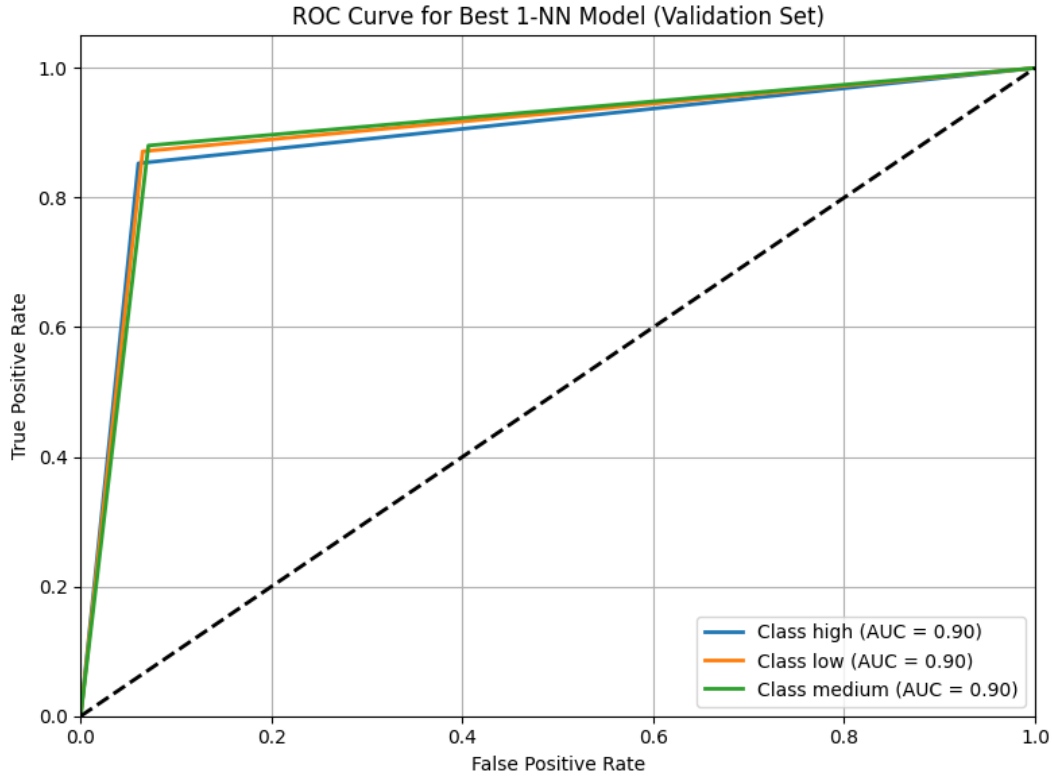


Figure 7: ROC curve for K-NN classifier with  $k=1$ .

The ROC curve for the K-NN classifier with  $k=1$  is shown in Figure 7. The AUC for the K-NN classifier with  $k=1$  is approximately 0.9 for all height categories. This indicates that the K-NN classifier with  $k=1$  is a good classifier, as it has a high true positive rate and a low false positive rate.

The ROC curve itself is not very informative for one-neighbour K-NN, as it can only produce three points. The lack of points is due to the fact that the model outputs "hard" probabilities (only 0 or 1) for each class, rather than "soft" probabilities, that would be possible for a higher number of neighbours.

## 5.4 Confusion Matrix

We can also use the same K-NN model to plot the confusion matrix using the validation data that was not used in the training or testing of the model. The confusion matrix is a table that is often used to describe the performance of a classification model on a set of data for which the true values are known. The confusion matrix shows the number of correct and incorrect predictions made by the model, broken down by class.

The confusion matrix shows that the diagonal elements to be the largest, which indicates that the K-NN classifier with  $k=1$  is able to correctly classify the majority of the data points.

The off-diagonal elements are much smaller, which indicates that the K-NN classifier with  $k=1$  is not making many incorrect predictions. This is a good sign, as it indicates that the K-NN classifier with  $k=1$  is a good classifier for this dataset.

The validation set has a significant imbalance in the number of data points in each category; similar to the raw training data, the medium category has the most data points, followed by the low and high categories. However, SMOTE was not applied to validation because artificially generating data points would

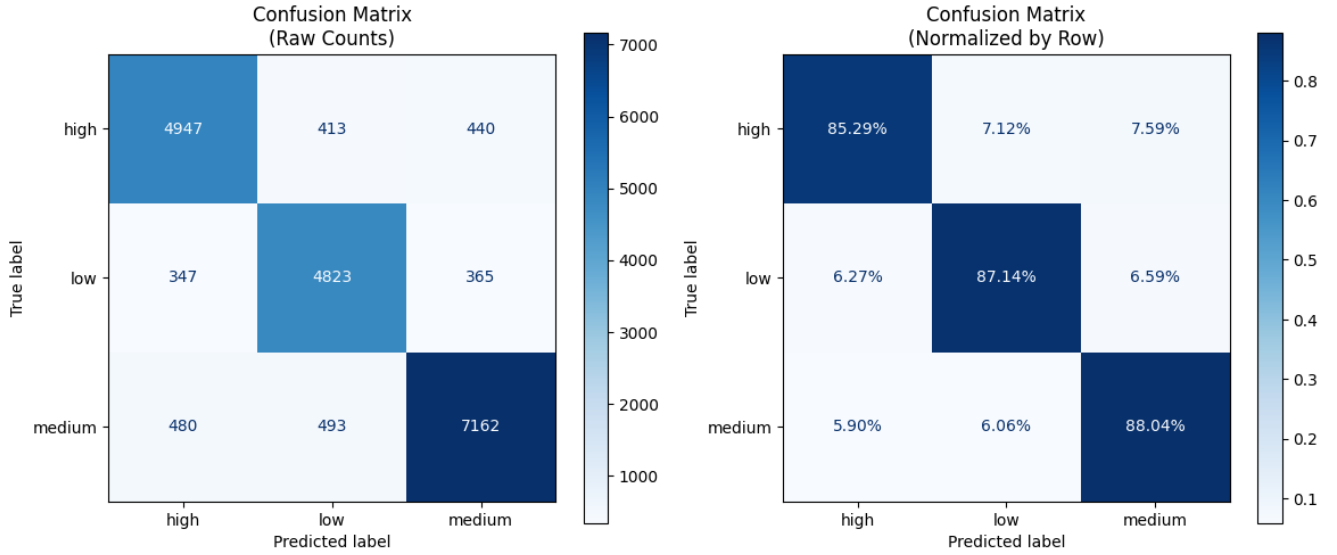


Figure 8: Confusion matrix for K-NN classifier with  $k=1$ . (Left) Confusion matrix against the raw validation set. (Right) Confusion matrix normalized by row.

give misleading performance metrics. Instead, the confusion matrix is normalized by row, which allows us to see the proportion of correct and incorrect predictions for each category. The normalized confusion matrix shows that the K-NN classifier with  $k=1$  is able to correctly classify the majority of the data points in each category.

## 5.5 Precision, Recall, and F1 score

The precision, recall, and F1 score are three metrics that are commonly used to evaluate the performance of a classification model. The precision is the number of true positive predictions divided by the total number of positive predictions. The recall is the number of true positive predictions divided by the total number of actual positive instances. The F1 score is the harmonic mean of precision and recall.

Table 1: Precision, Recall, and F1 Score for each class (Validation Set)

Class	Precision	Recall	F1-score
high	0.86	0.85	0.85
low	0.84	0.87	0.86
medium	0.90	0.88	0.89

All three classes have a high precision, recall, and F1 score, which indicates that the K-NN classifier with  $k=1$  is able to correctly classify the majority of the data points in each category. Additionally, the scores are relatively balanced across all classes, which suggests that the K-NN classifier with  $k=1$  is not biased towards any particular class, despite the imbalanced validation data.

## 6 Random Forest Classifier

The Random Forest classifier is an ensemble learning method that constructs a multitude of decision trees during training and outputs the mode of the classes (classification) or mean prediction (regression) of the individual trees. It is particularly effective for high-dimensional datasets and can handle both categorical and continuous variables. The Random Forest algorithm is robust to overfitting, especially when the number of trees in the forest is large. In this section, we will explore the performance of the Random Forest classifier on the TTSWING dataset.

Random Forest builds many decision trees using random subsets of the data and features, and aggregates their predictions. The quality of each split in a tree is measured using criteria such as Gini impurity (for classification), Entropy (information gain), or Mean Squared Error (for regression).

To keep this report within reasonable constraints, we will only explore the Gini impurity, and explore other parameters within the Random Forest classifier. The Gini impurity is a measure of how often a randomly chosen element from the set would be incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset.

### 6.1 Accuracy of Random Forest Classifier

The Random Forest classifier was trained on the balanced and scaled training data, and then tested on the testing data. Various values of the number of trees in the forest (10, 50, 100, 200, 500, 800, and 1400), and the maximum depth of the trees (5, 10, 20, 30, 40, and no limit) were explored. The accuracy of the Random Forest classifier was then calculated for each combination of parameters.

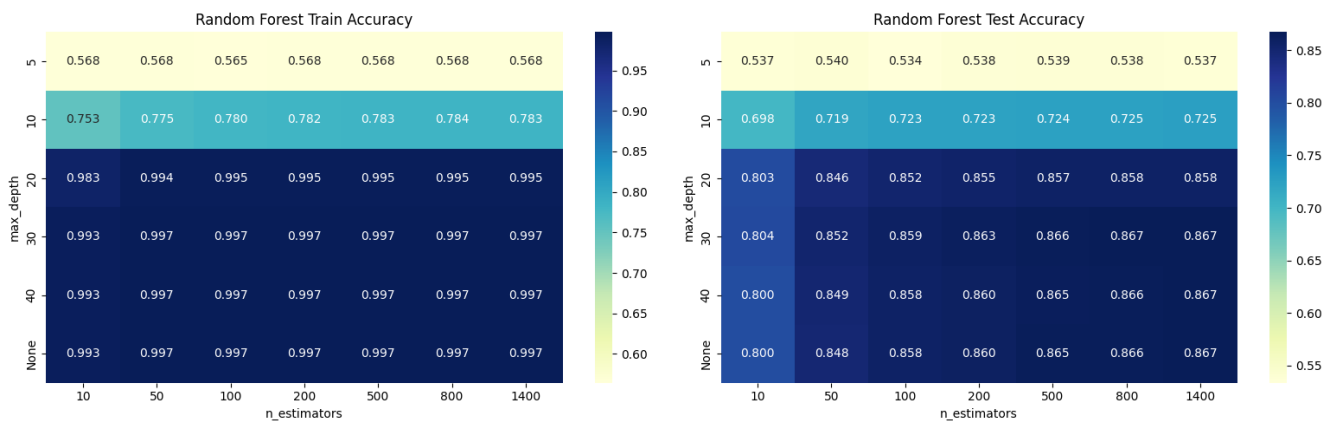


Figure 9: Accuracy of Random Forest classifier against the train and test data.

We can see from Figure 9 that the accuracy of the Random Forest classifier increases as the number of trees in the forest increases, and as the maximum depth of the trees increases. This is expected, as more trees and deeper trees allow for more complex decision boundaries to be learned. However, there is a point of diminishing returns, where increasing the number of trees or depth does not significantly improve accuracy.

As a caveat, the accuracy of the Random Forest classifier is not guaranteed to increase as the maximum depth of the trees increases. This is because deeper trees are more likely to overfit the data, and hence the accuracy may decrease as the maximum depth increases. This is particularly evident in all trees that were explored - the accuracy of the Random Forest classifier peaks at a maximum depth of either 20 or 30,

then decreases slightly as the maximum depth increases. This suggests that the Random Forest classifier is able to learn complex decision boundaries without overfitting the data, but that there is a limit to how complex the decision boundaries can be before overfitting occurs.

### 6.1.1 Computational time after dimensionality reduction

A big drawback with increasing the number of trees and the maximum depth of the trees is the computational time. The computational time of the Random Forest classifier is significantly increased as the number of trees and maximum depth increases. This is because each tree in the forest must be trained on the data, and each tree must be traversed to make a prediction. This can be mitigated by using dimensionality reduction techniques such as PCA, which reduces the number of features in the data, and hence reduces the computational time.

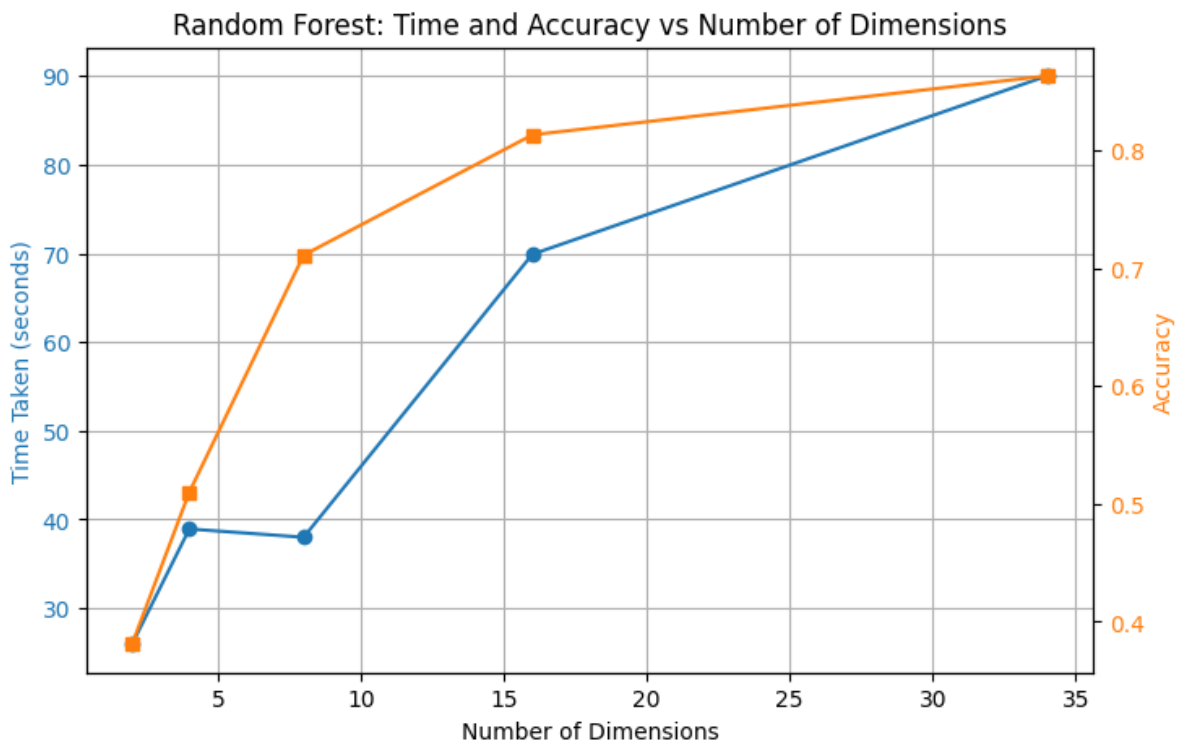


Figure 10: Accuracy of Random Forest classifier on PCA-reduced data. The model parameters were set to 200 trees and a maximum depth of 30.

When the data is reduced to 16-dimensions using PCA, the computational time of the Random Forest classifier significantly decreases. However, the accuracy of the Random Forest classifier slightly decreases as well. This is expected, as the PCA-reduced data loses some information content, and hence the accuracy of the Random Forest classifier decreases. However, the decrease in accuracy is not significant, and the computational time is significantly reduced. This is a strong argument for using PCA to reduce computational cost, without significantly sacrificing accuracy. However, since the time taken to train the Random Forest classifier on the original data is within reason, and the accuracy is slightly higher, we will use the original data for the rest of the analysis.

## 6.2 K-Fold Cross-validation

Increasing the number of trees does not increase risk of overfitting, however, increasing the maximum depth of the trees does increase the risk of overfitting. We can take advantage of this by deciding the best number of trees, since we are not limited by the number of trees with regards to overfitting. The accuracy of the Random Forest classifier only marginally increases after 100 trees, and the computational time increases exponentially after 100 trees. Therefore, we will use 100 trees and vary the maximum depth of the trees to investigate the effect of overfitting on the accuracy of the Random Forest classifier.

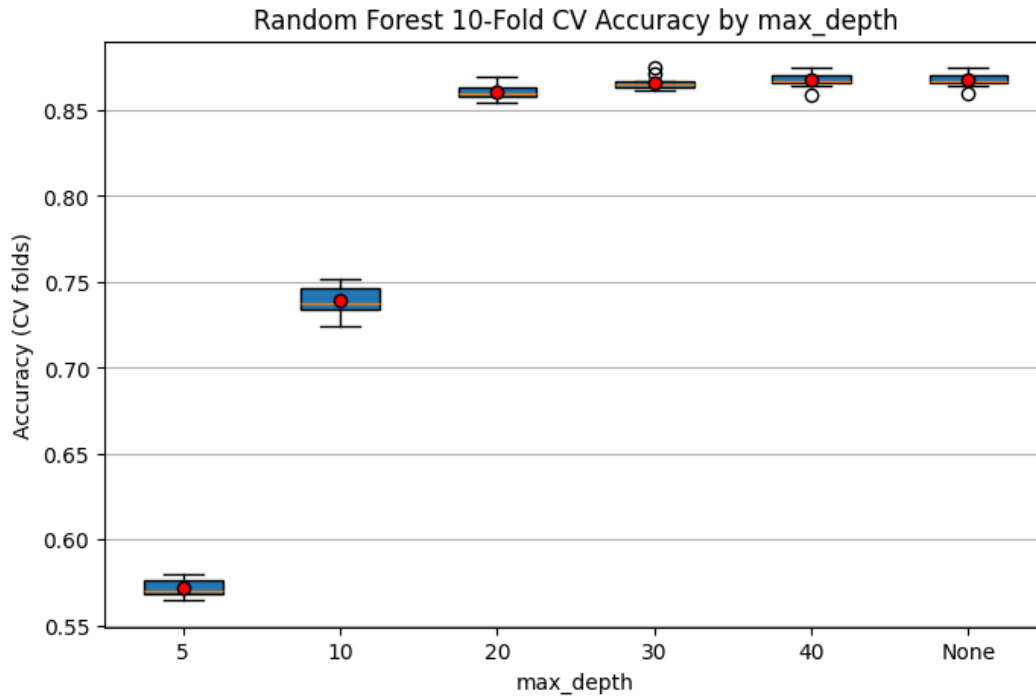


Figure 11: Accuracy of Random Forest classifier with K-Fold Cross-validation. Each box and whisker plots shows the spread of the accuracy against the test data across the 10 folds. The red line shows the mean accuracy across the 10 folds.

The distribution of the accuracy shows a low variance across all maximum depths, which suggests that the Random Forest classifier is not overfitting the data. The maximum depth of None (unlimited) achieved the highest mean accuracy, while still maintaining a low variance, therefore we will use this model as the standard model for the Random Forest classifier.

Note that the means of the accuracies for maximum depths larger than 20 are very similar. Additionally, the mean accuracy for 40 is within the IQR of the maximum depth of None. It is fair to assume the true mean across the maximum depths larger than 20 are insignificantly different, and so the choice of maximum depth greater than 20 is arbitrary, in terms of accuracy. Despite the fact that the computational time is significantly increased as the maximum depth increases, we will use the maximum depth of None, because the time taken is within reason, and the accuracy, for this dataset, is the highest. The Random Forest classifier with 100 trees and unlimited depth achieved a mean accuracy of approximately 86.75%, which is slightly higher than the K-NN classifier with  $k=1$ , which achieved a mean accuracy of approximately 86.26%. This suggests that the Random Forest classifier is a better classifier for this dataset than the K-NN classifier, and that the Random Forest classifier is able to learn more complex decision boundaries than the K-NN classifier.

### 6.3 ROC curve

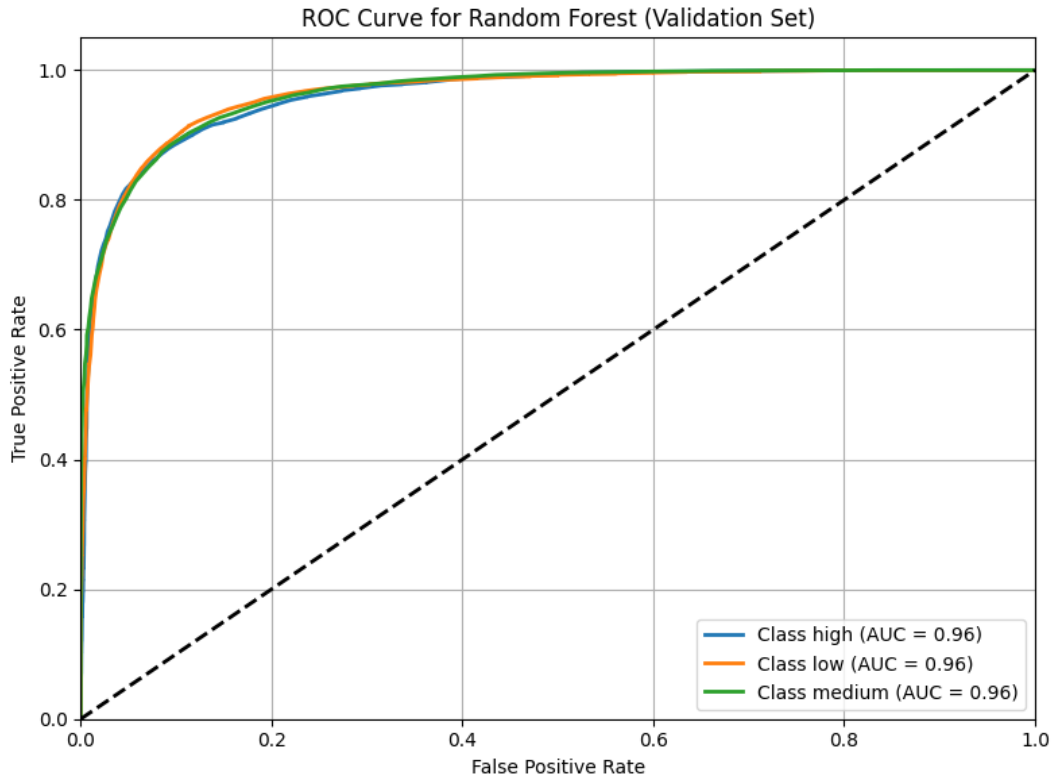


Figure 12: ROC curve for Random Forest classifier with 100 trees and maximum depth of None.

The ROC curve for the Random Forest classifier with 100 trees and unlimited depth is shown in Figure 12. The area under the curve (AUC) is approximately 0.96 for all height categories, indicating strong discriminative performance with a high true positive rate and a low false positive rate across classes. Notably, the Random Forest achieves a higher AUC than the K-NN classifier with ( $k=1$ ), suggesting superior classification performance on this dataset.

This improvement is consistent with the strengths of ensemble methods like Random Forest, which aggregate the predictions of multiple decision trees to reduce variance and improve generalization. Additionally, Random Forests are less prone to overfitting compared to single decision trees, particularly when a sufficient number of trees are used.



## 6.4 Confusion matrix

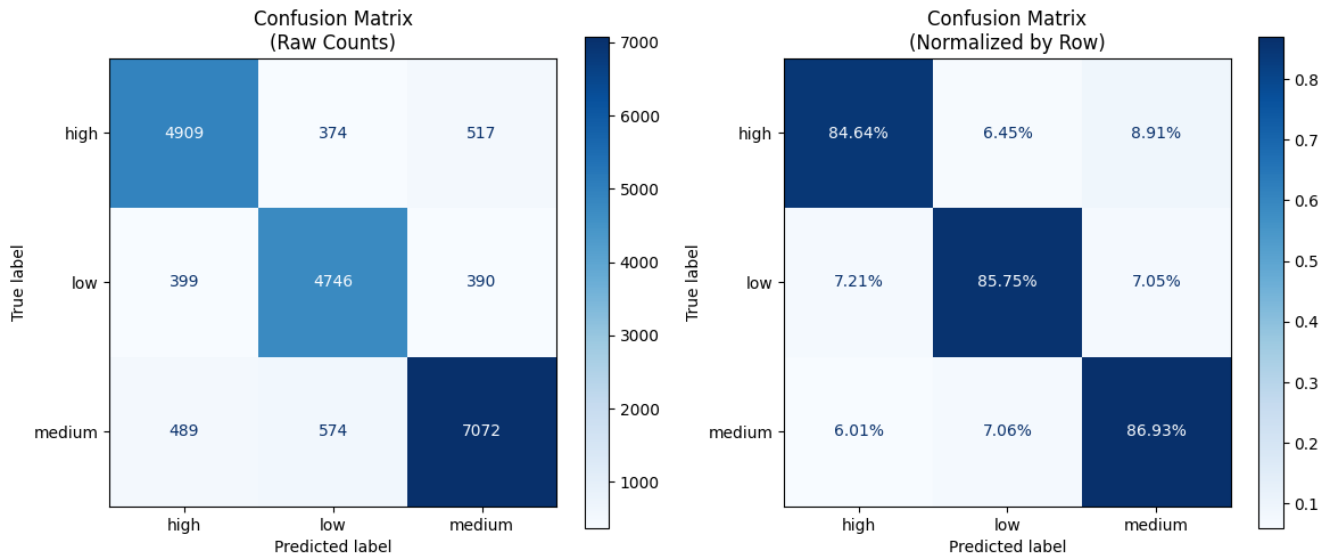


Figure 13: Confusion matrix for Random Forest classifier with 100 trees and maximum depth of None. (Left) Confusion matrix against the raw validation set. (Right) Confusion matrix normalized by row.

Similar to the K-NN classifier, the diagonal elements are the largest, which indicates that the Random Forest classifier with 100 trees and unlimited depth is able to correctly classify the majority of the data points.

The diagonal elements for the Random Forest classifier are slightly lower than the K-NN classifier, which indicates that the chosen Random Forest classifier is only slightly less accurate than the K-NN classifier.

## 6.5 Precision, Recall, and F1 score

Class	Precision	Recall	F1-score
high	0.85	0.85	0.85
low	0.83	0.86	0.85
medium	0.89	0.87	0.88

Table 2: Precision, recall, and F1-score for each class on the validation set.

The precision, recall, and F1 score are very similar to that of the K-NN classifier and are all high for all three classes, which indicates that the Random Forest classifier with 100 trees and unlimited depth is able to correctly classify the majority of the data points in each category. Additionally, the scores are relatively balanced across all classes, which suggests that the Random Forest classifier with 100 trees and unlimited depth is not biased towards any particular class.

## 7 Deep Neural Networks

The Deep Neural Network (DNN) is a powerful machine learning algorithm that can be used for both classification and regression tasks. The DNN algorithm works by using multiple layers of neurons to learn complex decision boundaries. For this section, we will explore two-hidden layer DNNs, which are a type of feedforward neural network. The model is trained using the Cross Entropy loss function, with a L2 regularization term to prevent overfitting. The model is trained using the Adam optimizer, which is a popular optimization algorithm for training neural networks. The model is trained on the scaled training data, and then tested on the testing data.

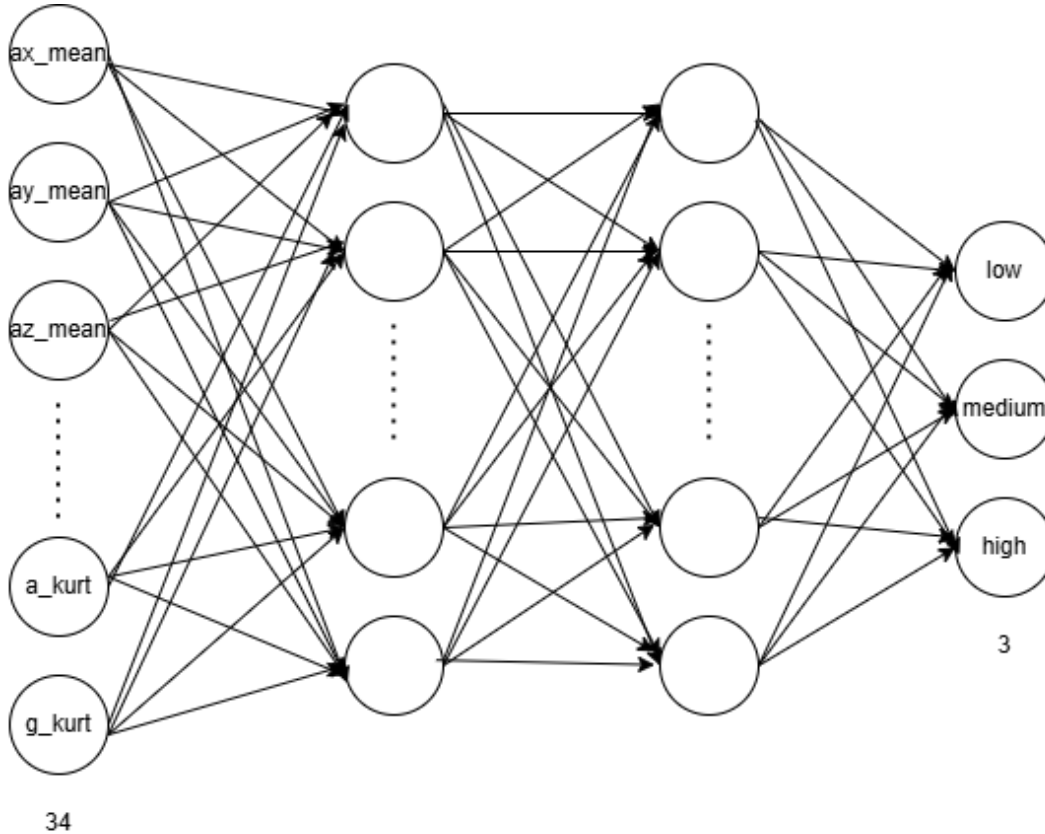


Figure 14: Architecture of the DNN with two hidden layers. The input layer has 34 nodes. There are two hidden layers with a variable number of nodes. The output layer has 3 nodes, corresponding to the three height categories.

### 7.1 Training

Initial assessment of the DNN were performed using a 128 nodes in each hidden layers. It was found that the DNN can achieved a accuracy on-par with the previous models when trained using a batch size of 16, learning rate of 0.001, and 100 epochs. However, it is unclear if the model is overfitting, and so we will use 5-fold cross-validation to investigate the effect of the different hidden layer sizes on the accuracy of the DNN.

### 7.1.1 Effect of hidden layer sizes

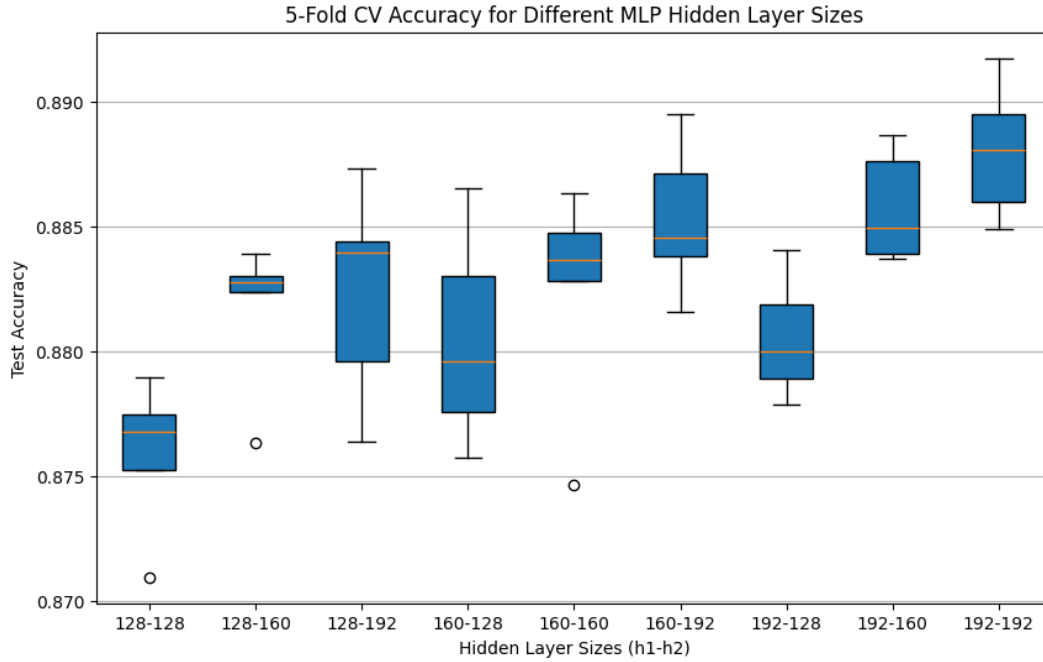


Figure 15: Different hidden layer sizes (h1-h2) with 5-fold cross-validation. Each box and whisker plots shows the spread of the accuracy against the test data across the 5 folds. The red line shows the mean accuracy across the 5 folds.

There appears to be a trend for increasing the size of the hidden layers to increase the accuracy of the DNN. Notably, the biggest increase in accuracies occurs when the second hidden layer is increased, for the same size of the first hidden layer. This suggests that the second hidden layer is more important than the first hidden layer, and that increasing the size of the second hidden layer will increase the accuracy of the DNN.

The DNN with 192-192 hidden layers has the highest mean accuracy, and has a low variance, that is comparable to the other models. This suggests that the DNN with 192-192 hidden layers is a good choice for this dataset, and that the DNN is not overfitting the data, or at least not overfitting the data more than the other models.

### 7.1.2 Training 192-192 DNN model

The best model structure from the 5-fold cross-validation is the DNN with 192-192 hidden layers, which has a mean accuracy of approximately 88.8%. We use this model as the standard model for the DNN, and train it on the full training data. The model is trained using a batch size of 16 and an initial learning rate of 0.001, and 100 epochs. The model was paused every 20 epochs to fine-tune the learning rate, and the model was trained for a total 116 epochs. After the learning rate was lowered, 100 epochs were trained, there was no significant increase in the train and test accuracy, so the model was stopped at 116 epochs. The learning rate was lowered during the training process to prevent overshooting the minimum of the loss function, while still training the model at a reasonable speed. The train/test curve are shown in Figure 16.

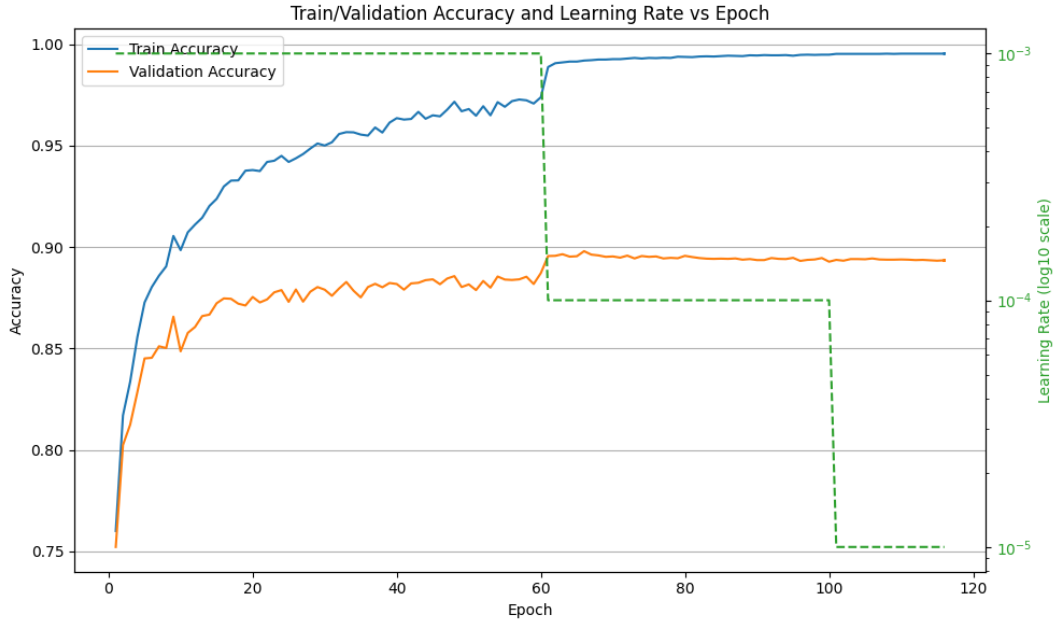


Figure 16: Train/test curve for the DNN with 192-192 hidden layers. The blue line shows the training accuracy, and the orange line shows the testing accuracy. The learning rate was lowered whenever the testing accuracy did not increase for 20 epochs. The model was trained for a total of 116 epochs.

## 7.2 ROC Curve

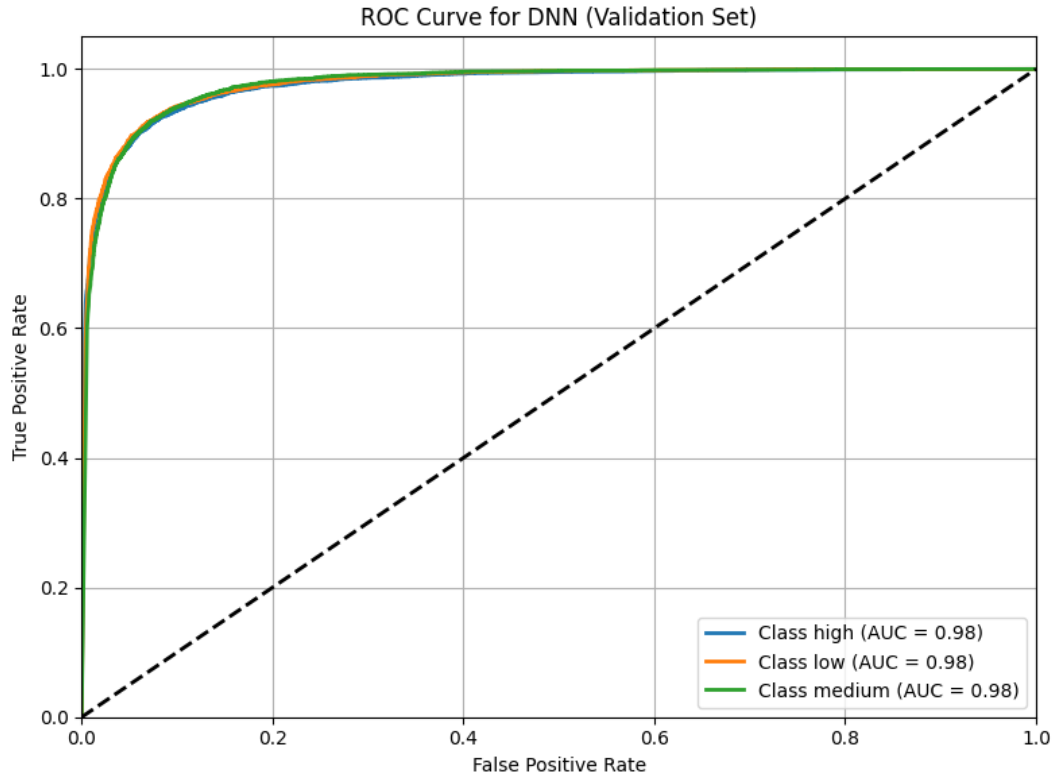


Figure 17: ROC curve for DNN with 192-192 hidden layers.

The area under the curve (AUC) is approximately 0.98 for all height categories, indicating strong discriminative performance with a high true positive rate and a low false positive rate across classes. The AUC is slightly higher than that of the Random Forest classifier, which suggests that the DNN is a better classifier for this dataset. This is expected, as the DNN is able to learn more complex decision boundaries than the Random Forest classifier, due to the multiple layers of neurons.

### 7.3 Confusion Matrix

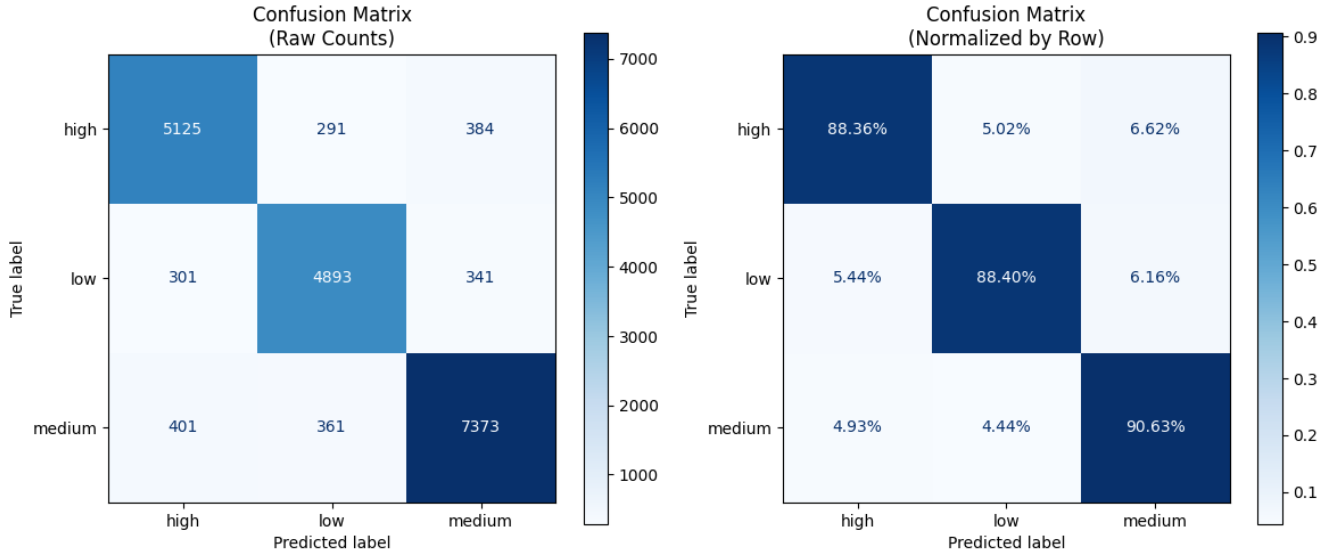


Figure 18: Confusion matrix for DNN with 192-192 hidden layers. (Left) Confusion matrix against the raw validation set. (Right) Confusion matrix normalized by row.

Similar to the K-NN and Random Forest classifiers, the diagonal elements are the largest, which indicates that the DNN with 192-192 hidden layers is able to correctly classify the majority of the data points. The diagonal elements for the DNN are slightly higher for all three classes, compared to the K-NN and Random Forest classifiers, which indicates that the DNN with 192-192 hidden layers is able to correctly classify more data points than the other two classifiers. This corroborates with the higher AUC of the DNN, and suggests that the DNN is a better classifier for this dataset.

### 7.4 Precision, Recall, and F1 score

Class	Precision	Recall	F1-score
high	0.88	0.88	0.88
low	0.88	0.88	0.88
medium	0.91	0.91	0.91

Table 3: Precision, Recall, and F1-score for each class on the validation set using the trained DNN model.

The precision, recall, and F1 scores are all high for all three classes, which indicates that the DNN with 192-192 hidden layers is able to correctly classify the majority of the data points in each category.

Additionally, the scores are relatively balanced across all classes, which suggests that the DNN with 192-192 hidden layers is not biased towards any particular class.

## 8 Conclusion

This report presented an analysis of the TTSWING dataset using three machine learning classifiers - K-Nearest Neighbours (K-NN), Random Forest, and Deep Neural Networks (DNN). The performance of each model was assessed using various metrics, including accuracy, ROC curve, confusion matrix, and precision, recall, and F1 score. Additional considerations such as data imbalance, dimensionality reduction, and model overfitting were also addressed.

The K-NN classifier achieved a maximum test accuracy of approximately 86.26% with one neighbour, and a ROC AUC of approximately 0.9 for all height categories. Dimensionality reduction using PCA was explored, but it decreased the accuracy of the K-NN classifier, suggesting that the curse of dimensionality was not a significant problem for this dataset. K-Fold cross-validation confirmed that the K-NN classifier with one neighbour did not overfit the data, as the variance in accuracy across folds was low.

The Random Forest classifier achieved a maximum test accuracy of approximately 86.75% with 100 trees and unlimited depth, and a ROC AUC of approximately 0.96 for all height categories - both metrics score higher than the K-classifier. Hyperparameter tuning was explored by varying the number of trees and maximum depth, with the best performance observed at 100 trees and unlimited depth. K-Fold cross-validation confirmed that the Random Forest classifier did not overfit the data, as the variance in accuracy across folds was low. PCA was also applied to the data, which significantly reduced the computational time of the Random Forest classifier, but did not significantly decrease the accuracy.

A two hidden-layer DNN outperformed both the K-NN and Random Forest classifiers, achieving a maximum test accuracy of approximately 88.8% with a ROC AUC of approximately 0.98 for all height categories. The DNN also demonstrated low variance in accuracy across folds during K-Fold cross-validation, indicating that it did not overfit the data. A hyperparameter search was conducted to find the optimal hidden layer sizes, with the best performance observed at 192-192 hidden layers.

While all three classifiers demonstrated strong performance at classifying the height of table tennis players based on their swing data, the DNN emerged as the most effective model in terms of predictive power, and balance across metrics. Nonetheless, trade-offs in training time and interpretability should be considered when these three models in practice. Here is a list of the key findings:

- K-NN classifier achieved a maximum test accuracy of 86.26% with one neighbour, and a ROC AUC of 0.9.
- Random Forest classifier achieved a maximum test accuracy of 86.75% with 100 trees and unlimited depth, and a ROC AUC of 0.96.
- DNN classifier achieved a maximum test accuracy of 88.8% with two hidden layers (192-192), and a ROC AUC of 0.98.
- Dimensionality reduction using PCA did not significantly improve the performance of the K-NN or Random Forest classifiers, but reduced computational time.

- K-Fold cross-validation confirmed that all three classifiers did not overfit the data, as the variance in accuracy across folds was low.

Future improvement could involve experimenting with different DNN architectures and incorporations of different structures such as convolutional layers, which could further enhance the model's performance. More extensive hyperparameter tuning for the Random Forest classifier, such as exploring different criteria for splitting nodes and specifying the minimum samples per leaf, should be explored. The K-NN classifier could be improved by exploring different distance metrics, such as Manhattan or Minkowski distance.

## 9 Appendix

### 9.1 All Project Code

Listing 1: All project code

```

1      #!/usr/bin/env python
2      # coding: utf-8
3
4      # ### Imports
5      #
6
7      # In[22]:
8
9
10     import pandas as pd
11     import numpy as np
12     from sklearn.preprocessing import StandardScaler
13     from sklearn.decomposition import PCA
14     from sklearn.cluster import KMeans
15     from sklearn.metrics import silhouette_score
16     import matplotlib.pyplot as plt
17     import seaborn as sns
18     from sklearn.model_selection import train_test_split
19     from sklearn.neighbors import KNeighborsClassifier
20     from sklearn.metrics import accuracy_score
21     from imblearn.over_sampling import SMOTE
22     from sklearn.preprocessing import LabelEncoder
23     import time
24     import torch
25     import torch.nn as nn
26     import torch.optim as optim
27     from sklearn.model_selection import StratifiedKFold
28     from sklearn.preprocessing import LabelEncoder
29     from sklearn.metrics import accuracy_score
30
31
32     # ### Data Exploration
33
34     # In[3]:
35
36
37     # Reload the CSV file

```

```

38 file_path = "assignTTSWING.csv"
39 df = pd.read_csv(file_path)
40
41
42 # In[4]:
43
44
45 # Filter only feature columns (exclude metadata and IDs)
46 feature_columns = [col for col in df.columns if col.startswith(('a_', 'g_', 'ax_',
47 'ay_', 'az_', 'gx_', 'gy_', 'gz_'))]
48 features = df[feature_columns]
49
50 # Assign the labels
51 testmode_labels = df['testmode'] if 'testmode' in df.columns else None
52 gender_labels = df['gender'] if 'gender' in df.columns else None
53 age_labels = df['age'] if 'age' in df.columns else None
54 playYears_labels = df['playYears'] if 'playYears' in df.columns else None
55 height_labels = df['height'] if 'height' in df.columns else None
56 weight_labels = df['weight'] if 'weight' in df.columns else None
57 handedness_labels = df['handedness'] if 'handedness' in df.columns else None
58 holdRacketHanded_labels = df['holdRacketHanded'] if 'holdRacketHanded' in df.
59 columns else None
60
61 # Combine all labels into a single DataFrame
62 labels_df = pd.DataFrame({
63     'testmode': testmode_labels,
64     'gender' : gender_labels,
65     'age' : age_labels,
66     'playYears' : playYears_labels,
67     'height' : height_labels,
68     'weight' : weight_labels,
69     'handedness' : handedness_labels,
70     'holdRacketHanded' : holdRacketHanded_labels
71 })
72
73 # In[5]:
74
75 # Convert all labels to string type
76 labels_df = labels_df.astype(str)
77
78
79 # In[6]:
80
81
82 # Count how many of each label we have
83 label_counts = labels_df.apply(pd.Series.value_counts).fillna(0).astype(int)
84
85 # Print the label counts
86 print("Label counts:")
87 print(label_counts)
88
89
90 # In[7]:
91

```



```

92
93 # Reveal the rows that are ??? for age, playYears, height, weight
94 missing_labels = labels_df[labels_df.isin(['???']).any(axis=1)]
95
96 # Print the rows with missing labels
97 print("\nRows with missing labels:")
98 print(missing_labels)
99
100 # Create a list of all the rows that are ??? for age, playYears, height, weight
101 missing_rows = missing_labels.index.tolist()
102
103 # Delete the rows with missing labels from the labels and features DataFrames
104 labels_df_cleaned = labels_df.drop(missing_rows)
105 features_cleaned = features.drop(missing_rows)
106
107
108 # In[8]:
109
110
111 # Recheck that the invalid rows are gone
112 # Count how many of each label we have
113 label_counts = labels_df_cleaned.apply(pd.Series.value_counts)
114
115 # Print the label counts
116 print("Label counts:")
117 print(label_counts)
118
119
120 # In[9]:
121
122
123 # Finally assign the X_train, X_test, X_val, y_train, y_test, y_val
124 # We will use 60% of the data for training, 20% for testing, and 20% for validation
125 # The y_train will be the height inside of the labels_df
126 X_train, X_temp, y_train, y_temp = train_test_split(features_cleaned,
127     labels_df_cleaned['height'], test_size=0.4, random_state=1)
128 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5,
129     random_state=1)
130
131
132 # In[10]:
133
134 # Check the sizes of the splits
135 print(f"X_train size: {X_train.shape}")
136 print(f"X_val size: {X_val.shape}")
137 print(f"X_test size: {X_test.shape}")
138
139 # Check the sizes of the splits
140 print(f"y_train size: {y_train.shape}")
141 print(f"y_val size: {y_val.shape}")
142 print(f"y_test size: {y_test.shape}")
143
144 # Check the ratios of the splits
145 print(f"X_train ratio: {X_train.shape[0] / features_cleaned.shape[0]}")
146 print(f"X_val ratio: {X_val.shape[0] / features_cleaned.shape[0]}")

```

```

146 print(f"X_test ratio: {X_test.shape[0] / features_cleaned.shape[0]}")
147
148 # Check the ratios of the splits
149 print(f"y_train ratio: {y_train.shape[0] / labels_df_cleaned.shape[0]}")
150 print(f"y_val ratio: {y_val.shape[0] / labels_df_cleaned.shape[0]}")
151 print(f"y_test ratio: {y_test.shape[0] / labels_df_cleaned.shape[0]}")
152
153
154 # ##### SMOTE
155
156 # In[11]:
157
158
159 # Check how many 'low', 'medium', and 'high' labels we have in the y_train
160 y_train_counts = y_train.value_counts()
161
162 # Print the y_train counts with percentage
163 y_train_counts_percentage = y_train_counts / y_train_counts.sum() * 100
164 print("\nY_train counts with percentage:")
165 print(y_train_counts_percentage)
166
167 print("\nY_train counts:")
168 print(y_train_counts)
169
170
171 # In[12]:
172
173
174 # Apply SMOTE only to training data
175 smote = SMOTE(random_state=42)
176 X_train_balanced, y_train_balanced = smote.fit_resample(X_train, y_train)
177
178 # Keep validation and test sets unchanged
179 # X_val, y_val remain as they are
180 # X_test, y_test remain as they are
181
182 # Print class distributions
183 print("Original training set distribution:", pd.Series(y_train).value_counts())
184 print("Balanced training set distribution:", pd.Series(y_train_balanced).
185       value_counts())
186 print("Validation set distribution:", pd.Series(y_val).value_counts())
187 print("Test set distribution:", pd.Series(y_test).value_counts())
188
189 # ##### STANDARDIZE THE FEATURES
190
191 # In[13]:
192
193
194 # Standardize the features
195 scaler = StandardScaler()
196 X_train_scaled = scaler.fit_transform(X_train_balanced)
197 X_val_scaled = scaler.transform(X_val)
198 X_test_scaled = scaler.transform(X_test)
199
200

```

```

201 # In[ ]:
202
203
204 import matplotlib.pyplot as plt
205
206 # Plot histograms for all features before standardization
207 features_cleaned.hist(bins=30, figsize=(18, 12), layout=(int(np.ceil(len(
    features_cleaned.columns)/6)), 6))
208 plt.suptitle('Feature Distributions Before Standardization', fontsize=16)
209 plt.tight_layout(rect=[0, 0, 1, 0.97])
210 plt.show()
211
212
213 # In[ ]:
214
215
216 import seaborn as sns
217 import matplotlib.pyplot as plt
218
219 corr = features_cleaned.corr()
220 plt.figure(figsize=(12, 10))
221 sns.heatmap(corr, cmap='coolwarm', center=0)
222 plt.title('Feature Correlation Matrix')
223 plt.show()
224
225
226 # ### PCA
227
228 # In[14]:
229
230
231 # Compute the PCA via the covariance matrix method
232
233 # Standardize the features
234 mean = features.mean()
235 std = features.std()
236 features_standardized = (features - mean) / std
237
238 # Compute the covariance matrix
239 cov_matrix = features_standardized.cov()
240
241 # Compute the eigenvalues and eigenvectors
242 eigenvalues, eigenvectors = np.linalg.eig(cov_matrix)
243
244 # Sort the eigenvalues and eigenvectors
245 sorted_indices = np.argsort(eigenvalues)[::-1]
246 sorted_eigenvalues = eigenvalues[sorted_indices]
247 sorted_eigenvectors = eigenvectors[:, sorted_indices]
248
249 # Perform dimensionality reduction by multiplying the data by the top two
    eigenvectors
250 pca_features = np.dot(features_standardized, sorted_eigenvectors[:, :2])
251 pca_features_df = pd.DataFrame(pca_features, columns=['PC1', 'PC2'])
252
253 # Plot the PCA results
254 plt.figure(figsize=(10, 6))

```

```

255 sns.scatterplot(data=pca_features_df, x='PC1', y='PC2', hue=testmode_labels,
    palette='Set1', alpha=0.7)
256 plt.title('PCA of Features')
257 plt.xlabel('Principal Component 1')
258 plt.ylabel('Principal Component 2')
259 plt.legend(title='Test Mode')
260 plt.grid()
261 plt.savefig('pca_features.png')
262 plt.show()
263
264
265
266 # In[15]:
267
268
269 # Again but with feature changed into ndarray
270 features_np = features.values
271
272 # Standardize the features
273 mean = np.mean(features_np, axis=0)
274 std_dev = np.std(features_np, axis=0)
275 features_standardized_np = (features_np - mean) / std_dev
276
277 # Compute the covariance matrix
278 cov_matrix_np = np.cov(features_standardized_np.T)
279
280 # Compute the eigenvalues and eigenvectors
281 eigenvalues_np, eigenvectors_np = np.linalg.eig(cov_matrix_np)
282 sorted_indices = np.argsort(eigenvalues_np)[::-1]
283 eigenvalues = eigenvalues_np[sorted_indices]
284 eigenvectors = eigenvectors_np[:, sorted_indices]
285
286 # Perform dimensionality reduction by multiplying the data by the top two
    eigenvectors
287 reduced_features = np.dot(features_standardized_np, eigenvectors[:, :2])
288
289
290
291 # In[16]:
292
293
294 # Make a Scree graph for all eigenvalues
295 plt.figure(figsize=(10, 6))
296 plt.plot(range(1, len(eigenvalues) + 1), eigenvalues, marker='o', linestyle='--')
297 plt.title('Scree Plot of Eigenvalues')
298 plt.xlabel('Principal Component')
299 plt.ylabel('Eigenvalue')
300 plt.xticks(range(1, len(eigenvalues) + 1))
301 plt.grid()
302 plt.savefig('scree_plot.png')
303 plt.show()
304
305 # Make another Scree graph for all eigenvalues but the y-axis is the percentage of
    variance explained
306 plt.figure(figsize=(10, 6))
307 explained_variance = eigenvalues / np.sum(eigenvalues) * 100

```

```

308 plt.plot(range(1, len(explained_variance) + 1), explained_variance, marker='o',
          linestyle='--')
309 plt.title('Scree Plot of Explained Variance')
310 plt.xlabel('Principal Component')
311 plt.ylabel('Explained Variance (%)')
312 plt.xticks(range(1, len(explained_variance) + 1))
313 plt.grid()
314 plt.show()
315
316 # Make a cumulative explained variance plot
317 plt.figure(figsize=(10, 6))
318 cumulative_explained_variance = np.cumsum(explained_variance)
319 plt.plot(range(1, len(cumulative_explained_variance) + 1),
          cumulative_explained_variance, marker='o', linestyle='--')
320 plt.title('Cumulative Explained Variance')
321 plt.xlabel('Principal Component')
322 plt.ylabel('Cumulative Explained Variance (%)')
323 plt.xticks(range(1, len(cumulative_explained_variance) + 1))
324 plt.grid()
325 plt.show()
326
327 # Check that the percentages add up to 100%
328 total_variance = np.sum(explained_variance)
329 print(f"Total variance explained by all components: {total_variance:.2f}%")
330
331
332 # In[17]:
333
334
335 # Reduce the features to 16 dimensions using PCA
336 reduced_features_16 = np.dot(features_standardized_np, eigenvectors[:, :16])
337
338 # Plot the correlation matrix of the reduced 16D features
339 plt.figure(figsize=(12, 10))
340 sns.heatmap(np.corrcoef(reduced_features_16.T), annot=True, fmt=".2f", cmap='
          coolwarm', square=True, cbar_kws={"shrink": .8})
341 plt.title('Correlation Matrix of Reduced 16D Features')
342 plt.show()
343
344
345 # In[18]:
346
347
348 # Standardize features
349 scaler = StandardScaler()
350 X_scaled = scaler.fit_transform(features)
351
352 # Apply PCA to reduce dimensions to 2D for visualization
353 pca = PCA(n_components=2)
354 X_pca = pca.fit_transform(X_scaled)
355
356
357 # ### K-NN ALGORITHM
358
359 # ##### K-NN WITH ORIGINAL DATASET
360

```

```

361 # In[19]:
362
363
364 def run_knn_from_split(X_train, X_test, y_train, y_test, n_neighbors=5, class_names
= None, plot_confusion=False, verbose=False, label_name=""):
365     # Fit KNN
366     knn = KNeighborsClassifier(n_neighbors=n_neighbors)
367     knn.fit(X_train, y_train)
368     y_pred = knn.predict(X_test)
369
370     # Accuracy and loss
371     accuracy = accuracy_score(y_test, y_pred)
372     loss = 1 - accuracy
373
374     if verbose:
375         print(f"\nK-NN Results for label: {label_name}")
376         print(classification_report(y_test, y_pred, target_names=class_names))
377         print(f"Accuracy: {accuracy:.4f}")
378         print(f"Loss: {loss:.4f}")
379
380     if plot_confusion:
381         from sklearn.metrics import confusion_matrix
382         cm = confusion_matrix(y_test, y_pred)
383         plt.figure(figsize=(6, 5))
384         sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
385                     xticklabels=class_names, yticklabels=class_names)
386         plt.xlabel('Predicted')
387         plt.ylabel('True')
388         plt.title(f'Confusion Matrix: {label_name}')
389         plt.tight_layout()
390         plt.show()
391
392     return accuracy, loss, knn
393
394 # # Example usage:
395 # from sklearn.model_selection import train_test_split
396 # from sklearn.preprocessing import LabelEncoder
397
398 # label_of_interest = 'height'
399 # y = labels_df[label_of_interest].values
400 # le = LabelEncoder()
401 # y_encoded = le.fit_transform(y)
402 # class_names = le.classes_
403
404 # # Split once
405 # X_train, X_test, y_train, y_test = train_test_split(
406 #     X_scaled, y_encoded, test_size=0.2, random_state=42, stratify=y_encoded
407 # )
408
409 # # Run KNN for k=1
410 # accuracy, loss, _ = run_knn_from_split(X_train, X_test, y_train, y_test,
411 #     n_neighbors=1, class_names=class_names, label_name=label_of_interest)
412
413 # In[20]:
414

```

```

415
416 n_neighbors_list = [2**i for i in range(13)]
417 results = []
418
419 label_of_interest = 'height' # any label you want
420
421 # Encode both training and test labels
422 le = LabelEncoder()
423 y_train_encoded = le.fit_transform(y_train_balanced) # Fit on training data
424 y_test_encoded = le.transform(y_test) # Transform test data using same encoder
425 class_names = le.classes_
426
427 for n in n_neighbors_list:
428     try:
429         accuracy, loss, knn = run_knn_from_split(
430             X_train_scaled,
431             X_test_scaled, # Use scaled test data
432             y_train_encoded, # Use encoded training labels
433             y_test_encoded, # Use encoded test labels
434             n_neighbors=n,
435             class_names=class_names,
436             label_name=label_of_interest
437         )
438         results.append({'n_neighbors': n, 'accuracy': accuracy, 'loss': loss, '
439                         knn_model': knn})
440         print(f"n_neighbors={n}: accuracy={accuracy:.4f}, loss={loss:.4f}")
441     except ValueError as e:
442         print(f"n_neighbors={n}: Error - {e}")
443
444 # Plot accuracy vs n_neighbors
445 plt.figure(figsize=(8, 5))
446 plt.plot([r['n_neighbors'] for r in results], [r['accuracy'] for r in results],
447          marker='o')
448 plt.xscale('log', base=2)
449 plt.xlabel('n_neighbors')
450 plt.ylabel('Accuracy')
451 plt.title(f'K-NN Accuracy vs n_neighbors ({label_of_interest})')
452 plt.grid(True)
453 plt.show()
454
455 # In[28]:
456
457 # Plot accuracy against the test AND training data
458 n_neighbors_list = [2**i for i in range(13)]
459 results = []
460
461 label_of_interest = 'height' # any label you want
462
463 # Encode both training and test labels
464 le = LabelEncoder()
465 y_train_encoded = le.fit_transform(y_train_balanced) # Fit on training data
466 y_test_encoded = le.transform(y_test) # Transform test data using same encoder
467 class_names = le.classes_
468

```

```

469 for n in n_neighbors_list:
470     try:
471         knn = KNeighborsClassifier(n_neighbors=n)
472         knn.fit(X_train_scaled, y_train_encoded)
473         y_train_pred = knn.predict(X_train_scaled)
474         y_test_pred = knn.predict(X_test_scaled)
475         train_acc = accuracy_score(y_train_encoded, y_train_pred)
476         test_acc = accuracy_score(y_test_encoded, y_test_pred)
477         results.append({'n_neighbors': n, 'train_accuracy': train_acc, '
                        test_accuracy': test_acc, 'knn_model': knn})
478         print(f"n_neighbors={n}: train_acc={train_acc:.4f}, test_acc={test_acc:.4f}
              ")
479     except ValueError as e:
480         print(f"n_neighbors={n}: Error - {e}")
481
482 # Plot training and test accuracy vs n_neighbors
483 plt.figure(figsize=(8, 5))
484 plt.plot([r['n_neighbors'] for r in results], [r['train_accuracy'] for r in results
         ], marker='o', label='Train Accuracy')
485 plt.plot([r['n_neighbors'] for r in results], [r['test_accuracy'] for r in results
         ], marker='s', label='Test Accuracy')
486 plt.xscale('log', base=2)
487 plt.xlabel('n_neighbors')
488 plt.ylabel('Accuracy')
489 plt.title(f'K-NN Train/Test Accuracy vs n_neighbors ({label_of_interest})')
490 plt.legend()
491 plt.grid(True)
492 plt.show()
493
494
495 # ##### K-NN CROSS-VALIDATION
496 #
497 #
498
499 # In[22]:
500
501
502 from sklearn.model_selection import StratifiedKFold, cross_val_score
503 from imblearn.pipeline import Pipeline
504 from imblearn.over_sampling import SMOTE
505 from sklearn.neighbors import KNeighborsClassifier
506
507 # Combine X_train and X_test for X input into K-Fold
508 X_kfold = np.vstack((X_train, X_test))
509 y_fold = np.hstack((y_train, y_test))
510
511 # Scale X_kfold
512 X_kfold_scaled = scaler.fit_transform(X_kfold)
513
514 k_values = [2**i for i in range(6)]
515 cv = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
516
517 # Store the results
518 results_kfold = []
519 all_scores = [] # Store all fold scores for each k
520

```



```

521 for k in k_values:
522     pipeline = Pipeline([
523         ('smote', SMOTE(random_state=42)),
524         ('knn', KNeighborsClassifier(n_neighbors=k))
525     ])
526     scores = cross_val_score(pipeline, X_kfold_scaled, y_fold, cv=cv, scoring='
accuracy')
527     print(f"k={k}: mean accuracy={scores.mean():.4f} (+/- {scores.std():.4f})")
528     results_kfold.append({
529         'k': k,
530         'mean_accuracy': scores.mean(),
531         'std_accuracy': scores.std(),
532         'min_accuracy': scores.min(),
533         'max_accuracy': scores.max()
534     })
535     all_scores.append(scores) # Save the array of scores for boxplot
536
537
538 # In[25]:
539
540
541 plt.figure(figsize=(8, 5))
542 plt.boxplot(
543     all_scores,
544     positions=range(len(k_values)),
545     widths=0.5,
546     patch_artist=True,
547     boxprops=dict(facecolor='lightblue', color='blue'),
548     whiskerprops=dict(color='black'),
549     capprops=dict(color='black'),
550     medianprops=dict(color='red')
551 )
552
553 plt.xticks(
554     range(len(k_values)),
555     [f"$2^{i}$" for i in range(len(k_values))]
556 )
557 plt.xlabel('k (number of neighbors)')
558 plt.ylabel('Accuracy (CV folds)')
559 plt.title('K-MN Accuracy Distribution with K-Fold Cross-Validation')
560 plt.grid(True, axis='y')
561 plt.show()
562
563
564 # #### K-NN WITH PCA REDUCTION
565
566 # In[27]:
567
568
569 # Compare K-NN performance with different PCA dimensions
570 pca_dimensions = [16, 14, 12, 10, 8, 6, 4, 2]
571 n_neighbors_list = [2**i for i in range(13)]
572 results_multi_pca = {d: [] for d in pca_dimensions}
573 computation_times = {}
574
575 # First run KNN on original data

```

```

576 start_time = time.time()
577 results_original = []
578
579 for n in n_neighbors_list:
580     try:
581         accuracy, loss, knn = run_knn_from_split(
582             X_train_scaled,
583             X_test_scaled,
584             y_train_encoded,
585             y_test_encoded,
586             n_neighbors=n,
587             class_names=class_names,
588             label_name="Original (33D)"
589         )
590         results_original.append({'n_neighbors': n, 'accuracy': accuracy, 'loss':
591                                 loss})
592         print(f"Original data - n_neighbors={n}: accuracy={accuracy:.4f}")
593     except ValueError as e:
594         print(f"Original data - n_neighbors={n}: Error - {e}")
595
596 computation_times[33] = time.time() - start_time
597
598 # Prepare all PCA transformations first
599 pca_data = {}
600 for n_components in pca_dimensions:
601     pca = PCA(n_components=n_components)
602     X_train_pca = pca.fit_transform(X_train_scaled)
603     X_test_pca = pca.transform(X_test_scaled)
604     pca_data[n_components] = (X_train_pca, X_test_pca)
605
606 # Run KNN for each PCA dimension
607 for n_components in pca_dimensions:
608     X_train_pca, X_test_pca = pca_data[n_components]
609
610     # Start timing only K-NN part
611     start_time = time.time()
612
613     for n in n_neighbors_list:
614         try:
615             accuracy, loss, knn = run_knn_from_split(
616                 X_train_pca,
617                 X_test_pca,
618                 y_train_encoded,
619                 y_test_encoded,
620                 n_neighbors=n,
621                 class_names=class_names,
622                 label_name=f"PCA-{n_components}D"
623             )
624             results_multi_pca[n_components].append(
625                 {'n_neighbors': n, 'accuracy': accuracy, 'loss': loss}
626             )
627             print(f"PCA-{n_components}D - n_neighbors={n}: accuracy={accuracy:.4f}")
628         except ValueError as e:
629             print(f"PCA-{n_components}D - n_neighbors={n}: Error - {e}")

```

```

630     # Record only K-NN computation time
631     computation_times[n_components] = time.time() - start_time
632
633     # Plot results for this PCA dimension
634     plt.plot([r['n_neighbors'] for r in results_multi_pca[n_components]],
635             [r['accuracy'] for r in results_multi_pca[n_components]],
636             marker='o', label=f'PCA ({n_components}D)', linestyle='--')
637
638 plt.xscale('log', base=2)
639 plt.xlabel('k (number of neighbors)')
640 plt.ylabel('Accuracy')
641 plt.title(f'Comparison of Different PCA Dimensions')
642 plt.grid(True)
643 plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
644 plt.tight_layout()
645 plt.show()
646
647 # Plot computation times
648 plt.figure(figsize=(10, 6))
649 sorted_dimensions = sorted(computation_times.keys(), reverse=True) # Sort
    dimensions in descending order
650 times = [computation_times[d] for d in sorted_dimensions]
651 plt.plot(sorted_dimensions, times, marker='o')
652 plt.xlabel('Number of Dimensions')
653 plt.ylabel('Computation Time (seconds)')
654 plt.title('Computation Time vs Dimensionality')
655 plt.grid(True)
656 plt.show()
657
658 # Print computation times
659 print("\nComputation times:")
660 print(f"Original (33D): {computation_times[33]:.2f} seconds")
661 for dim in pca_dimensions:
662     print(f"PCA-{dim}D: {computation_times[dim]:.2f} seconds")
663
664
665 # ##### K-NN ROC curve
666
667 # In[17]:
668
669
670 from sklearn.model_selection import StratifiedKFold
671 from imblearn.pipeline import Pipeline
672 from imblearn.over_sampling import SMOTE
673 from sklearn.neighbors import KNeighborsClassifier
674 from sklearn.metrics import roc_curve, auc
675 from sklearn.preprocessing import LabelBinarizer
676 import numpy as np
677 import matplotlib.pyplot as plt
678
679 # Prepare data for cross-validation (already in your code)
680 X_kfold = np.vstack((X_train, X_test))
681 y_fold = np.hstack((y_train, y_test))
682 X_kfold_scaled = scaler.fit_transform(X_kfold)
683 cv = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
684

```

```

685 # Store models and scores for k=1
686 models = []
687 scores = []
688
689 for train_idx, test_idx in cv.split(X_kfold_scaled, y_fold):
690     X_tr, X_te = X_kfold_scaled[train_idx], X_kfold_scaled[test_idx]
691     y_tr, y_te = y_fold[train_idx], y_fold[test_idx]
692     # Pipeline with SMOTE and 1-NN
693     pipeline = Pipeline([
694         ('smote', SMOTE(random_state=42)),
695         ('knn', KNeighborsClassifier(n_neighbors=1))
696     ])
697     pipeline.fit(X_tr, y_tr)
698     acc = pipeline.score(X_te, y_te)
699     models.append(pipeline)
700     scores.append(acc)
701
702 # Find the best model (highest accuracy on its fold)
703 best_idx = np.argmax(scores)
704 best_model = models[best_idx]
705 print(f"Best fold accuracy: {scores[best_idx]:.4f}")
706
707 # Use the best model to predict probabilities on the validation set
708 y_val_proba = best_model.predict_proba(X_val_scaled)
709
710 # Binarize y_val for multiclass ROC
711 lb = LabelBinarizer()
712 y_val_binarized = lb.fit_transform(y_val)
713 if y_val_binarized.shape[1] == 1:
714     y_val_binarized = np.hstack([1 - y_val_binarized, y_val_binarized])
715
716 # Plot ROC curve for each class
717 plt.figure(figsize=(8, 6))
718 for i in range(y_val_binarized.shape[1]):
719     fpr, tpr, _ = roc_curve(y_val_binarized[:, i], y_val_proba[:, i])
720     roc_auc = auc(fpr, tpr)
721     plt.plot(fpr, tpr, lw=2, label=f'Class {lb.classes_[i]} (AUC = {roc_auc:.2f})')
722
723 plt.plot([0, 1], [0, 1], 'k--', lw=2)
724 plt.xlim([0.0, 1.0])
725 plt.ylim([0.0, 1.05])
726 plt.xlabel('False Positive Rate')
727 plt.ylabel('True Positive Rate')
728 plt.title('ROC Curve for Best 1-NN Model (Validation Set)')
729 plt.legend(loc="lower right")
730 plt.grid(True)
731 plt.tight_layout()
732 plt.show()
733
734
735 # #### K-NN Confusion Matrix
736
737 # In[18]:
738
739
740 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

```

```

741 # Predict on the validation set
742 y_val_pred = best_model.predict(X_val_scaled)
743
744 # If you used label encoding, ensure y_val_pred and y_val are comparable
745 # If y_val is not encoded, encode it with the same encoder as used in training
746 # Example:
747 # y_val_encoded = le.transform(y_val)
748
749 # Compute confusion matrix
750 cm = confusion_matrix(y_val, y_val_pred, labels=best_model.classes_)
751
752 # Plot confusion matrix
753 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=best_model.
754                               classes_)
755 plt.figure(figsize=(6, 5))
756 disp.plot(cmap='Blues', values_format='d')
757 plt.title('Confusion Matrix: Best 1-NN Model (Validation Set)')
758 plt.tight_layout()
759 plt.show()
760
761 # In[19]:
762
763 # Normalized confusion matrix
764 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
765 import numpy as np
766
767 # Predict on the validation set
768 y_val_pred = best_model.predict(X_val_scaled)
769
770 # Compute confusion matrices (raw counts and normalized)
771 cm = confusion_matrix(y_val, y_val_pred, labels=best_model.classes_)
772 cm_normalized = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
773
774 # Create side-by-side plots
775 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 5))
776
777 # Plot raw counts
778 disp1 = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=best_model.
779                                classes_)
780 disp1.plot(ax=ax1, cmap='Blues', values_format='d')
781 ax1.set_title('Confusion Matrix\n(Raw Counts)')
782
783 # Plot normalized values
784 disp2 = ConfusionMatrixDisplay(confusion_matrix=cm_normalized, display_labels=
785                                best_model.classes_)
786 disp2.plot(ax=ax2, cmap='Blues', values_format='.2%')
787 ax2.set_title('Confusion Matrix\n(Normalized by Row)')
788
789 plt.tight_layout()
790 plt.show()
791
792 # Print classification rates for each class
793 print("\nClassification rates per class:")

```

```

794 for i, class_name in enumerate(best_model.classes_):
795     print(f"{class_name}: {cm_normalized[i][i]:.2%} correct")
796
797
798 # In[20]:
799
800
801 from sklearn.metrics import classification_report
802
803 # Predict on the validation set
804 y_val_pred = best_model.predict(X_val_scaled)
805
806 # Print precision, recall, and F1 score for each class
807 print("Precision, Recall, and F1 Score (Validation Set):\n")
808 print(classification_report(y_val, y_val_pred, target_names=best_model.classes_))
809
810
811 # ### RANDOM FOREST CLASSIFIER
812
813 # ##### Random Forest on original dataset
814
815 # In[23]:
816
817
818 import numpy as np
819 import pandas as pd
820 from sklearn.ensemble import RandomForestClassifier
821 from sklearn.metrics import accuracy_score
822 import matplotlib.pyplot as plt
823 import seaborn as sns
824
825 # Define parameter ranges
826 n_estimators_list = [10, 50, 100, 200, 500]
827 max_depth_list = [None, 5, 10, 20, 30, 40]
828
829 # Store results
830 results = []
831
832 for n_estimators in n_estimators_list:
833     for max_depth in max_depth_list:
834         rf = RandomForestClassifier(
835             n_estimators=n_estimators,
836             max_depth=max_depth,
837             criterion='gini',
838             random_state=42
839         )
840         rf.fit(X_train_scaled, y_train_balanced)
841         y_pred = rf.predict(X_test_scaled)
842         acc = accuracy_score(y_test, y_pred)
843         results.append({
844             'n_estimators': n_estimators,
845             'max_depth': max_depth if max_depth is not None else 'None',
846             'accuracy': acc
847         })
848     print(f"n_estimators={n_estimators}, max_depth={max_depth}: accuracy={acc:.4f}")

```

```

849
850 # Convert results to DataFrame for heatmap
851 df_results = pd.DataFrame(results)
852 pivot_table = df_results.pivot(index='max_depth', columns='n_estimators', values='
    accuracy')
853
854 # Plot heatmap
855 plt.figure(figsize=(8, 5))
856 sns.heatmap(pivot_table, annot=True, fmt=".3f", cmap="YlGnBu")
857 plt.title("Random Forest Accuracy\n(varying n_estimators and max_depth)")
858 plt.xlabel("n_estimators")
859 plt.ylabel("max_depth")
860 plt.show()
861
862
863 # In[24]:
864
865
866 # Continue with more results
867
868 # Define parameter ranges
869 n_estimators_list = [800, 1400]
870 max_depth_list = [None, 5, 10, 20, 30, 40]
871
872 # Store results
873 results_2 = []
874
875 for n_estimators in n_estimators_list:
876     for max_depth in max_depth_list:
877         rf = RandomForestClassifier(
878             n_estimators=n_estimators,
879             max_depth=max_depth,
880             criterion='gini',
881             random_state=42
882         )
883         rf.fit(X_train_scaled, y_train_balanced)
884         y_pred = rf.predict(X_test_scaled)
885         acc = accuracy_score(y_test, y_pred)
886         results_2.append({
887             'n_estimators': n_estimators,
888             'max_depth': max_depth if max_depth is not None else 'None',
889             'accuracy': acc
890         })
891         print(f"n_estimators={n_estimators}, max_depth={max_depth}: accuracy={acc
            :.4f}")
892
893 # Combine results with results_2 to make a combined results
894 combined_results = results + results_2
895
896 # Convert results to DataFrame for heatmap
897 df_results = pd.DataFrame(combined_results)
898 pivot_table = df_results.pivot(index='max_depth', columns='n_estimators', values='
    accuracy')
899
900 # Plot heatmap
901 plt.figure(figsize=(8, 5))

```

```

902 sns.heatmap(pivot_table, annot=True, fmt=".3f", cmap="YlGnBu")
903 plt.title("Random Forest Accuracy\n(varying n_estimators and max_depth)")
904 plt.xlabel("n_estimators")
905 plt.ylabel("max_depth")
906 plt.show()
907
908
909 # In[33]:
910
911
912 # SAME AS ABOVE BUT WITH TRAINING AND TEST ACCURACY
913
914 # Define parameter ranges
915 n_estimators_list = [10, 50, 100, 200, 500, 800, 1400]
916 max_depth_list = [None, 5, 10, 20, 30, 40]
917
918 # Store results
919 results = []
920
921 for n_estimators in n_estimators_list:
922     for max_depth in max_depth_list:
923         rf = RandomForestClassifier(
924             n_estimators=n_estimators,
925             max_depth=max_depth,
926             criterion='gini',
927             random_state=42
928         )
929         rf.fit(X_train_scaled, y_train_balanced)
930         y_train_pred = rf.predict(X_train_scaled)
931         y_test_pred = rf.predict(X_test_scaled)
932         train_acc = accuracy_score(y_train_balanced, y_train_pred)
933         test_acc = accuracy_score(y_test, y_test_pred)
934         results.append({
935             'n_estimators': n_estimators,
936             'max_depth': max_depth if max_depth is not None else 'None',
937             'train_accuracy': train_acc,
938             'test_accuracy': test_acc
939         })
940         print(f"n_estimators={n_estimators}, max_depth={max_depth}: train_acc={
941             train_acc:.4f}, test_acc={test_acc:.4f}")
942
943
944 # In[34]:
945
946
947 # Convert results to DataFrame for heatmap
948 df_results = pd.DataFrame(results)
949
950 # Pivot for heatmaps
951 pivot_train = df_results.pivot(index='max_depth', columns='n_estimators', values='
952     train_accuracy')
953 pivot_test = df_results.pivot(index='max_depth', columns='n_estimators', values='
954     test_accuracy')
955
956 # Plot heatmaps

```



```

955 plt.figure(figsize=(16, 5))
956 plt.subplot(1, 2, 1)
957 sns.heatmap(pivot_train, annot=True, fmt=".3f", cmap="YlGnBu")
958 plt.title("Random Forest Train Accuracy")
959 plt.xlabel("n_estimators")
960 plt.ylabel("max_depth")
961
962 plt.subplot(1, 2, 2)
963 sns.heatmap(pivot_test, annot=True, fmt=".3f", cmap="YlGnBu")
964 plt.title("Random Forest Test Accuracy")
965 plt.xlabel("n_estimators")
966 plt.ylabel("max_depth")
967
968 plt.tight_layout()
969 plt.show()
970
971
972 # #### Random Forest on PCA reduced dataset
973
974 # In[ ]:
975
976
977 # List of dimensions to test (including original)
978 dimensions = [X_train_scaled.shape[1], 16, 8, 4, 2]
979 times = []
980 accuracies = []
981
982 for dim in dimensions:
983     if dim == X_train_scaled.shape[1]:
984         # Use original data
985         X_train_dim = X_train_scaled
986         X_test_dim = X_test_scaled
987     else:
988         # PCA reduction
989         pca = PCA(n_components=dim, random_state=42)
990         X_train_dim = pca.fit_transform(X_train_scaled)
991         X_test_dim = pca.transform(X_test_scaled)
992
993     # Measure time
994     start = time.time()
995     rf = RandomForestClassifier(n_estimators=200, max_depth=30, criterion='gini',
996                               random_state=42)
997     rf.fit(X_train_dim, y_train_balanced)
998     y_pred = rf.predict(X_test_dim)
999     elapsed = time.time() - start
1000     acc = accuracy_score(y_test, y_pred)
1001
1002     times.append(elapsed)
1003     accuracies.append(acc)
1004     print(f"Dimensions: {dim}, Time: {elapsed:.2f}s, Accuracy: {acc:.4f}")
1005
1006 # # Plot time vs dimensions
1007 # plt.figure(figsize=(8, 5))
1008 # plt.plot(dimensions, times, marker='o')
1009 # plt.xlabel('Number of Dimensions')
1010 # plt.ylabel('Time Taken (seconds)')

```

```

1010 # plt.title('Random Forest: Time vs Number of Dimensions')
1011 # plt.grid(True)
1012 # plt.show()
1013
1014 # # Plot accuracy vs dimensions
1015 # plt.figure(figsize=(8, 5))
1016 # plt.plot(dimensions, accuracies, marker='o')
1017 # plt.xlabel('Number of Dimensions')
1018 # plt.ylabel('Accuracy')
1019 # plt.title('Random Forest: Accuracy vs Number of Dimensions')
1020 # plt.grid(True)
1021 # plt.show()
1022
1023
1024 # In[27]:
1025
1026
1027 plt.figure(figsize=(8, 5))
1028 ax1 = plt.gca()
1029 color1 = 'tab:blue'
1030 color2 = 'tab:orange'
1031
1032 # Plot time taken (left y-axis)
1033 ax1.plot(dimensions, times, marker='o', color=color1, label='Time Taken (s)')
1034 ax1.set_xlabel('Number of Dimensions')
1035 ax1.set_ylabel('Time Taken (seconds)', color=color1)
1036 ax1.tick_params(axis='y', labelcolor=color1)
1037
1038 # Create a second y-axis for accuracy
1039 ax2 = ax1.twinx()
1040 ax2.plot(dimensions, accuracies, marker='s', color=color2, label='Accuracy')
1041 ax2.set_ylabel('Accuracy', color=color2)
1042 ax2.tick_params(axis='y', labelcolor=color2)
1043
1044 plt.title('Random Forest: Time and Accuracy vs Number of Dimensions')
1045 ax1.grid(True)
1046 plt.show()
1047
1048
1049 # ##### Random Forest with cross-validation
1050
1051 # In[35]:
1052
1053
1054 from sklearn.ensemble import RandomForestClassifier
1055 from sklearn.model_selection import cross_val_score, StratifiedKFold
1056 import matplotlib.pyplot as plt
1057 import numpy as np
1058
1059 # Combine train and test sets
1060 X_kfold = np.vstack((X_train, X_test))
1061 y_fold = np.hstack((y_train, y_test))
1062
1063 max_depth_list = [5, 10, 20, 30, 40, None]
1064 cv = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
1065 all_scores = []

```

```

1066
1067 for max_depth in max_depth_list:
1068     rf = RandomForestClassifier(
1069         n_estimators=100,
1070         max_depth=max_depth,
1071         criterion='gini',
1072         random_state=42
1073     )
1074     scores = cross_val_score(rf, X_kfold, y_fold, cv=cv, scoring='accuracy')
1075     all_scores.append(scores)
1076     print(f"max_depth={max_depth}: mean={scores.mean():.4f}, std={scores.std():.4f}")
1077
1078 # Plot box-and-whisker plots with mean
1079 plt.figure(figsize=(8, 5))
1080 plt.boxplot(
1081     all_scores,
1082     labels=[str(md) for md in max_depth_list],
1083     patch_artist=True,
1084     showmeans=True,
1085     meanprops=dict(marker='o', markerfacecolor='red', markeredgecolor='black')
1086 )
1087 plt.xlabel('max_depth')
1088 plt.ylabel('Accuracy (CV folds)')
1089 plt.title('Random Forest 10-Fold CV Accuracy by max_depth')
1090 plt.grid(True, axis='y')
1091 plt.show()
1092
1093
1094 # ##### Random Forest ROC Curve
1095
1096 # In[37]:
1097
1098
1099 from sklearn.ensemble import RandomForestClassifier
1100 from sklearn.metrics import roc_curve, auc
1101 from sklearn.preprocessing import LabelBinarizer
1102 import matplotlib.pyplot as plt
1103 import numpy as np
1104
1105 # Train the best model
1106 rf = RandomForestClassifier(n_estimators=100, max_depth=None, random_state=42)
1107 rf.fit(X_train_scaled, y_train_balanced)
1108
1109 # Predict probabilities on validation set
1110 y_val_proba = rf.predict_proba(X_val_scaled)
1111
1112 # Binarize the validation labels for multiclass ROC
1113 lb = LabelBinarizer()
1114 y_val_binarized = lb.fit_transform(y_val)
1115 if y_val_binarized.shape[1] == 1:
1116     y_val_binarized = np.hstack([1 - y_val_binarized, y_val_binarized])
1117
1118 # Plot ROC curve for each class
1119 plt.figure(figsize=(8, 6))
1120 for i in range(y_val_binarized.shape[1]):

```

```

1121     fpr, tpr, _ = roc_curve(y_val_binarized[:, i], y_val_proba[:, i])
1122     roc_auc = auc(fpr, tpr)
1123     plt.plot(fpr, tpr, lw=2, label=f'Class {lb.classes_[i]} (AUC = {roc_auc:.2f})')
1124
1125 plt.plot([0, 1], [0, 1], 'k--', lw=2)
1126 plt.xlim([0.0, 1.0])
1127 plt.ylim([0.0, 1.05])
1128 plt.xlabel('False Positive Rate')
1129 plt.ylabel('True Positive Rate')
1130 plt.title('ROC Curve for Random Forest (Validation Set)')
1131 plt.legend(loc="lower right")
1132 plt.grid(True)
1133 plt.tight_layout()
1134 plt.show()
1135
1136
1137 # ##### Random Forest Confusion Matrix
1138
1139 # In[39]:
1140
1141
1142 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
1143 import numpy as np
1144
1145 # Predict on the validation set
1146 y_val_pred = rf.predict(X_val_scaled)
1147
1148 # Compute confusion matrices (raw counts and normalized)
1149 cm = confusion_matrix(y_val, y_val_pred, labels=rf.classes_)
1150 cm_normalized = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
1151
1152 # Create side-by-side plots
1153 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 5))
1154
1155 # Plot raw counts
1156 disp1 = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=rf.classes_)
1157 disp1.plot(ax=ax1, cmap='Blues', values_format='d')
1158 ax1.set_title('Confusion Matrix\n(Raw Counts)')
1159
1160 # Plot normalized values
1161 disp2 = ConfusionMatrixDisplay(confusion_matrix=cm_normalized, display_labels=rf.
    classes_)
1162 disp2.plot(ax=ax2, cmap='Blues', values_format='.2%')
1163 ax2.set_title('Confusion Matrix\n(Normalized by Row)')
1164
1165 plt.tight_layout()
1166 plt.show()
1167
1168 # Print classification rates for each class
1169 print("\nClassification rates per class:")
1170 for i, class_name in enumerate(rf.classes_):
1171     print(f"{class_name}: {cm_normalized[i][i]:.2%} correct")
1172
1173
1174 # ##### Random Forest Precision, Recall, F1
1175

```

```

1176 # In[40]:
1177
1178
1179 from sklearn.metrics import classification_report
1180
1181 # Predict on the validation set
1182 y_val_pred = rf.predict(X_val_scaled)
1183
1184 # Print precision, recall, and F1 score for each class
1185 print("Precision, Recall, and F1 Score (Validation Set):\n")
1186 print(classification_report(y_val, y_val_pred, target_names=rf.classes_))
1187
1188
1189 # ### Deep Neural Net
1190
1191 # ##### Training and testing two hidden layers with Cross Entropy (CE) and L2
1192     regularisation
1193
1194 # In[ ]:
1195
1196 import torch
1197 import torch.nn as nn
1198 import torch.optim as optim
1199 from sklearn.preprocessing import LabelEncoder
1200 from sklearn.metrics import accuracy_score
1201
1202 # Define the model
1203 class MLP(nn.Module):
1204     def __init__(self, input_dim, hidden1, hidden2, num_classes):
1205         super(MLP, self).__init__()
1206         self.fc1 = nn.Linear(input_dim, hidden1)
1207         self.fc2 = nn.Linear(hidden1, hidden2)
1208         self.fc3 = nn.Linear(hidden2, num_classes)
1209         self.relu = nn.ReLU()
1210     def forward(self, x):
1211         x = self.relu(self.fc1(x))
1212         x = self.relu(self.fc2(x))
1213         x = self.fc3(x)
1214         return x
1215
1216 # Set your variables
1217 hidden1 = 128 # number of nodes in first hidden layer
1218 hidden2 = 128 # number of nodes in second hidden layer
1219 lambda_l2 = 0.00 # L2 regularization strength
1220 epochs = 100
1221 batch_size = 16
1222 learning_rate = 0.001
1223 epochs_per_round = 100 # number of epochs to train before asking for user input
1224
1225 # Use GPU if available
1226 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
1227
1228 # Encode labels
1229 le = LabelEncoder()
1230 y_train_enc = le.fit_transform(y_train_balanced)

```

```

1231 y_test_enc = le.transform(y_test)
1232 num_classes = len(le.classes_)
1233
1234 # Convert data to torch tensors
1235 X_train_tensor = torch.tensor(X_train_scaled, dtype=torch.float32)
1236 y_train_tensor = torch.tensor(y_train_enc, dtype=torch.long)
1237 X_test_tensor = torch.tensor(X_test_scaled, dtype=torch.float32)
1238 y_test_tensor = torch.tensor(y_test_enc, dtype=torch.long)
1239
1240 # Dataset and DataLoader
1241 train_dataset = torch.utils.data.TensorDataset(X_train_tensor, y_train_tensor)
1242 train_loader = torch.utils.data.DataLoader(train_dataset, batch_size=batch_size,
1243                                           shuffle=True)
1244
1245 model = MLP(X_train_scaled.shape[1], hidden1, hidden2, num_classes).to(device)
1246
1247 # Loss and optimizer
1248 criterion = nn.CrossEntropyLoss()
1249 optimizer = optim.Adam(model.parameters(), lr=learning_rate, weight_decay=lambda_l2
1250 )
1251
1252 # Training loop
1253 total_epochs = 0
1254 while True:
1255     for epoch in range(epochs_per_round):
1256         model.train()
1257         for xb, yb in train_loader:
1258             xb, yb = xb.to(device), yb.to(device)
1259             optimizer.zero_grad()
1260             outputs = model(xb)
1261             loss = criterion(outputs, yb)
1262             loss.backward()
1263             optimizer.step()
1264         total_epochs += 1
1265         if (epoch+1) % 5 == 0 or epoch == 0:
1266             print(f"Epoch {total_epochs}, Loss: {loss.item():.4f}")
1267
1268     # Ask user if they want to continue
1269     cont = input("Continue training for another 30 epochs? (y/n): ").strip().lower()
1270     if cont != 'y':
1271         break
1272
1273     # Ask for new learning rate
1274     try:
1275         new_lr = float(input("Enter new learning rate for next 30 epochs (current:
1276                               {:.5f}): ".format(optimizer.param_groups[0]['lr'])))
1277         for param_group in optimizer.param_groups:
1278             param_group['lr'] = new_lr
1279     except Exception:
1280         print("Invalid input. Keeping previous learning rate.")
1281
1282 # Evaluate on test set
1283 model.eval()
1284 with torch.no_grad():
1285     X_test_tensor = X_test_tensor.to(device)
1286     y_test_tensor = y_test_tensor.to(device)

```

```

1283     outputs = model(X_test_tensor)
1284     _, predicted = torch.max(outputs, 1)
1285     test_acc = accuracy_score(y_test_tensor.cpu(), predicted.cpu())
1286     print(f"Test accuracy: {test_acc:.4f}")
1287
1288
1289 # ##### Deep neural net 5-fold cross-validation
1290
1291 # In[ ]:
1292
1293
1294 # --- Prepare data ---
1295 # Combine train and test sets for K-Fold
1296 X_kfold = np.vstack((X_train, X_test))
1297 y_fold = np.hstack((y_train, y_test))
1298 # Scale X_kfold
1299 X_kfold_scaled = scaler.fit_transform(X_kfold)
1300
1301 X = X_kfold_scaled
1302 le = LabelEncoder()
1303 y = le.fit_transform(y_fold)
1304 num_classes = len(le.classes_)
1305
1306 # --- Model definition ---
1307 class MLP(nn.Module):
1308     def __init__(self, input_dim, hidden1, hidden2, num_classes):
1309         super().__init__()
1310         self.fc1 = nn.Linear(input_dim, hidden1)
1311         self.fc2 = nn.Linear(hidden1, hidden2)
1312         self.fc3 = nn.Linear(hidden2, num_classes)
1313         self.relu = nn.ReLU()
1314     def forward(self, x):
1315         x = self.relu(self.fc1(x))
1316         x = self.relu(self.fc2(x))
1317         x = self.fc3(x)
1318         return x
1319
1320 # --- Hyperparameters ---
1321 hidden_sizes = [128, 160, 192]
1322 lambda_l2 = 0.00
1323 epochs = 100
1324 batch_size = 16
1325 learning_rate = 0.001
1326 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
1327
1328 # --- Cross-validation setup ---
1329 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
1330 results = {} # {(h1, h2): [fold accuracies]}
1331
1332 for h1 in hidden_sizes:
1333     for h2 in hidden_sizes:
1334         fold accuracies = []
1335         print(f"\nTraining MLP with hidden sizes: ({h1}, {h2})")
1336         for fold, (train_idx, val_idx) in enumerate(skf.split(X, y)):
1337             X_tr, X_val = X[train_idx], X[val_idx]
1338             y_tr, y_val = y[train_idx], y[val_idx]

```

```

1339
1340     # Convert to tensors
1341     X_tr_tensor = torch.tensor(X_tr, dtype=torch.float32)
1342     y_tr_tensor = torch.tensor(y_tr, dtype=torch.long)
1343     X_val_tensor = torch.tensor(X_val, dtype=torch.float32)
1344     y_val_tensor = torch.tensor(y_val, dtype=torch.long)
1345
1346     train_dataset = torch.utils.data.TensorDataset(X_tr_tensor, y_tr_tensor
1347 )
1348     train_loader = torch.utils.data.DataLoader(train_dataset, batch_size=
1349 batch_size, shuffle=True)
1350
1351     model = MLP(X.shape[1], h1, h2, num_classes).to(device)
1352     criterion = nn.CrossEntropyLoss()
1353     optimizer = optim.Adam(model.parameters(), lr=learning_rate,
1354 weight_decay=lambda_l2)
1355
1356     # Training loop
1357     for epoch in range(epochs):
1358         model.train()
1359         for xb, yb in train_loader:
1360             xb, yb = xb.to(device), yb.to(device)
1361             optimizer.zero_grad()
1362             outputs = model(xb)
1363             loss = criterion(outputs, yb)
1364             loss.backward()
1365             optimizer.step()
1366             if (epoch+1) % 10 == 0 or epoch == 0:
1367                 print(f" Fold {fold+1}, Epoch {epoch+1}/{epochs}, Loss: {loss.
1368 item():.4f}")
1369
1370     # Validation
1371     model.eval()
1372     with torch.no_grad():
1373         X_val_tensor = X_val_tensor.to(device)
1374         y_val_tensor = y_val_tensor.to(device)
1375         outputs = model(X_val_tensor)
1376         _, predicted = torch.max(outputs, 1)
1377         acc = accuracy_score(y_val_tensor.cpu(), predicted.cpu())
1378         fold_accuracies.append(acc)
1379         print(f" Fold {fold+1}: Accuracy = {acc:.4f}")
1380
1381     results[(h1, h2)] = fold_accuracies
1382
1383 # In[66]:
1384
1385 # --- Box and whisker plot ---
1386 plt.figure(figsize=(10, 6))
1387 labels = []
1388 data = []
1389 for (h1, h2), accs in results.items():
1390     labels.append(f"{h1}-{h2}")
1391     data.append(accs)

```



```

1391 plt.boxplot(data, labels=labels, patch_artist=True)
1392 plt.xlabel("Hidden Layer Sizes (h1-h2)")
1393 plt.ylabel("Test Accuracy")
1394 plt.title("5-Fold CV Accuracy for Different MLP Hidden Layer Sizes")
1395 plt.grid(True, axis='y')
1396 plt.show()
1397
1398
1399 # #### Training best model (192-192)
1400
1401 # In[ ]:
1402
1403
1404 # Hyperparameters
1405 input_dim = X_train_scaled.shape[1]
1406 hidden1 = 192
1407 hidden2 = 192
1408 epochs = 500
1409 learning_rate = 0.001
1410 lambda_l2 = 0.00
1411 patience = 1000
1412 batch_size = 16
1413 prompt_every = 20 # Number of epochs to train before prompting user
1414 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
1415
1416 # Model definition
1417 class MLP(nn.Module):
1418     def __init__(self, input_dim, hidden1, hidden2, num_classes):
1419         super().__init__()
1420         self.fc1 = nn.Linear(input_dim, hidden1)
1421         self.fc2 = nn.Linear(hidden1, hidden2)
1422         self.fc3 = nn.Linear(hidden2, num_classes)
1423         self.relu = nn.ReLU()
1424     def forward(self, x):
1425         x = self.relu(self.fc1(x))
1426         x = self.relu(self.fc2(x))
1427         x = self.fc3(x)
1428         return x
1429
1430 # Encode labels
1431 le = LabelEncoder()
1432 y_train_enc = le.fit_transform(y_train_balanced)
1433 y_val_enc = le.transform(y_val)
1434 num_classes = len(le.classes_)
1435
1436 # Convert data to torch tensors
1437 X_train_tensor = torch.tensor(X_train_scaled, dtype=torch.float32)
1438 y_train_tensor = torch.tensor(y_train_enc, dtype=torch.long)
1439 X_val_tensor = torch.tensor(X_val_scaled, dtype=torch.float32)
1440 y_val_tensor = torch.tensor(y_val_enc, dtype=torch.long)
1441
1442 # DataLoader
1443 train_dataset = torch.utils.data.TensorDataset(X_train_tensor, y_train_tensor)
1444 train_loader = torch.utils.data.DataLoader(train_dataset, batch_size=batch_size,
1445                                             shuffle=True)

```

```

1446 model = MLP(input_dim, hidden1, hidden2, num_classes).to(device)
1447 criterion = nn.CrossEntropyLoss()
1448 optimizer = optim.Adam(model.parameters(), lr=learning_rate, weight_decay=lambda_12
    )
1449
1450 # Early stopping variables
1451 best_val_acc = 0.0
1452 epochs_no_improve = 0
1453 best_model_state = None
1454
1455 train_accuracies = []
1456 val_accuracies = []
1457
1458 total_epochs = 0
1459 while total_epochs < epochs:
1460     for _ in range(prompt_every):
1461         if total_epochs >= epochs:
1462             break
1463         model.train()
1464         for xb, yb in train_loader:
1465             xb, yb = xb.to(device), yb.to(device)
1466             optimizer.zero_grad()
1467             outputs = model(xb)
1468             loss = criterion(outputs, yb)
1469             loss.backward()
1470             optimizer.step()
1471         # Training accuracy
1472         model.eval()
1473         with torch.no_grad():
1474             train_outputs = model(X_train_tensor.to(device))
1475             _, train_pred = torch.max(train_outputs, 1)
1476             train_acc = accuracy_score(y_train_tensor.cpu(), train_pred.cpu())
1477             train_accuracies.append(train_acc)
1478         # Validation accuracy
1479         val_outputs = model(X_val_tensor.to(device))
1480         _, val_pred = torch.max(val_outputs, 1)
1481         val_acc = accuracy_score(y_val_tensor.cpu(), val_pred.cpu())
1482         val_accuracies.append(val_acc)
1483         print(f"Epoch {total_epochs+1}, Loss: {loss.item():.4f}, Train Acc: {
            train_acc:.4f}, Val Acc: {val_acc:.4f}, Learning Rate: {optimizer.
            param_groups[0]['lr']:.5f}")
1484         # Early stopping check
1485         if val_acc > best_val_acc:
1486             best_val_acc = val_acc
1487             epochs_no_improve = 0
1488             best_model_state = model.state_dict()
1489         else:
1490             epochs_no_improve += 1
1491         if epochs_no_improve >= patience:
1492             print(f"Early stopping at epoch {total_epochs+1}. Best Val Acc: {
                best_val_acc:.4f}")
1493             break
1494         total_epochs += 1
1495     # Prompt user for learning rate change or to continue
1496     if total_epochs < epochs and epochs_no_improve < patience:

```

```

1497     cont = input(f"Continue training for another {prompt_every} epochs? (y/n):
1498         ").strip().lower()
1499     if cont != 'y':
1500         break
1501     try:
1502         new_lr = float(input(f"Enter new learning rate (current: {optimizer.
1503             param_groups[0]['lr']:.5f}): "))
1504         for param_group in optimizer.param_groups:
1505             param_group['lr'] = new_lr
1506     except Exception:
1507         print("Invalid input. Keeping previous learning rate.")
1508
1509 # In[33]:
1510
1511 import re
1512 import matplotlib.pyplot as plt
1513
1514 # Path to your log file
1515 log_path = "dnn_train_test_results.yaml"
1516
1517 # Prepare lists
1518 epochs = []
1519 train_accs = []
1520 val_accs = []
1521 lrs = []
1522
1523 # Regex to extract values
1524 epoch_re = re.compile(
1525     r"Epoch (\d+), Loss: [\d\.eE+-]+, Train Acc: ([\d\.eE+-]+), Val Acc: ([\d\.eE+-]+), Learning Rate: ([\d\.eE+-]+)"
1526 )
1527
1528 with open(log_path, "r") as f:
1529     for line in f:
1530         match = epoch_re.search(line)
1531         if match:
1532             epoch = int(match.group(1))
1533             train_acc = float(match.group(2))
1534             val_acc = float(match.group(3))
1535             lr = float(match.group(4))
1536             epochs.append(epoch)
1537             train_accs.append(train_acc)
1538             val_accs.append(val_acc)
1539             lrs.append(lr)
1540
1541 # Plot
1542 fig, ax1 = plt.subplots(figsize=(10, 6))
1543
1544 color1 = 'tab:blue'
1545 color2 = 'tab:orange'
1546 color3 = 'tab:green'
1547
1548 # Plot train/val accuracy

```

```

1550 ax1.plot(epochs, train_accs, label='Train Accuracy', color=color1)
1551 ax1.plot(epochs, val_accs, label='Validation Accuracy', color=color2)
1552 ax1.set_xlabel('Epoch')
1553 ax1.set_ylabel('Accuracy')
1554 ax1.legend(loc='upper left')
1555 ax1.grid(True, axis='y')
1556
1557 # Plot learning rate on right y-axis (log10 scale)
1558 ax2 = ax1.twinx()
1559 ax2.plot(epochs, lrs, label='Learning Rate', color=color3, linestyle='--')
1560 ax2.set_ylabel('Learning Rate (log10 scale)', color=color3)
1561 ax2.tick_params(axis='y', labelcolor=color3)
1562 ax2.set_yscale('log')
1563
1564 plt.title('Train/Validation Accuracy and Learning Rate vs Epoch')
1565 plt.tight_layout()
1566 plt.show()
1567
1568
1569 # In[30]:
1570
1571
1572 # Load best model
1573 if best_model_state is not None:
1574     model.load_state_dict(best_model_state)
1575
1576
1577 # ##### ROC Curve
1578
1579 # In[34]:
1580
1581
1582 from sklearn.metrics import roc_curve, auc
1583 from sklearn.preprocessing import LabelBinarizer
1584 import matplotlib.pyplot as plt
1585 import numpy as np
1586 import torch
1587
1588 # Ensure model is in eval mode and on correct device
1589 model.eval()
1590 X_val_tensor = torch.tensor(X_val_scaled, dtype=torch.float32).to(device)
1591
1592 # Get predicted probabilities from the DNN
1593 with torch.no_grad():
1594     logits = model(X_val_tensor)
1595     y_val_proba = torch.softmax(logits, dim=1).cpu().numpy()
1596
1597 # Binarize the validation labels for multiclass ROC
1598 lb = LabelBinarizer()
1599 y_val_binarized = lb.fit_transform(y_val)
1600 if y_val_binarized.shape[1] == 1:
1601     y_val_binarized = np.hstack([1 - y_val_binarized, y_val_binarized])
1602
1603 # Plot ROC curve for each class
1604 plt.figure(figsize=(8, 6))
1605 for i in range(y_val_binarized.shape[1]):

```

```

1606     fpr, tpr, _ = roc_curve(y_val_binarized[:, i], y_val_proba[:, i])
1607     roc_auc = auc(fpr, tpr)
1608     plt.plot(fpr, tpr, lw=2, label=f'Class {lb.classes_[i]} (AUC = {roc_auc:.2f})')
1609
1610 plt.plot([0, 1], [0, 1], 'k--', lw=2)
1611 plt.xlim([0.0, 1.0])
1612 plt.ylim([0.0, 1.05])
1613 plt.xlabel('False Positive Rate')
1614 plt.ylabel('True Positive Rate')
1615 plt.title('ROC Curve for DNN (Validation Set)')
1616 plt.legend(loc="lower right")
1617 plt.grid(True)
1618 plt.tight_layout()
1619 plt.show()
1620
1621
1622 # #### Confusion Matrix
1623
1624 # In[36]:
1625
1626
1627 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
1628 import numpy as np
1629
1630 # Predict on the validation set using the best DNN model
1631 if best_model_state is not None:
1632     model.load_state_dict(best_model_state)
1633 model.eval()
1634 with torch.no_grad():
1635     X_val_tensor = torch.tensor(X_val_scaled, dtype=torch.float32).to(device)
1636     logits = model(X_val_tensor)
1637     y_val_pred = torch.argmax(logits, dim=1).cpu().numpy()
1638
1639 # Use the same label encoder as during training for class names
1640 class_names = le.classes_
1641
1642 # Compute confusion matrices (raw counts and normalized)
1643 cm = confusion_matrix(y_val_enc, y_val_pred, labels=range(len(class_names)))
1644 cm_normalized = cm.astype('float') / cm.sum(axis=1, keepdims=True)
1645
1646 # Create side-by-side plots
1647 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 5))
1648
1649 # Plot raw counts
1650 disp1 = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=class_names)
1651 disp1.plot(ax=ax1, cmap='Blues', values_format='d')
1652 ax1.set_title('Confusion Matrix\n(Raw Counts)')
1653
1654 # Plot normalized values
1655 disp2 = ConfusionMatrixDisplay(confusion_matrix=cm_normalized, display_labels=
    class_names)
1656 disp2.plot(ax=ax2, cmap='Blues', values_format='.2%')
1657 ax2.set_title('Confusion Matrix\n(Normalized by Row)')
1658
1659 plt.tight_layout()
1660 plt.show()

```

```

1661
1662 # Print classification rates for each class
1663 print("\nClassification rates per class:")
1664 for i, class_name in enumerate(class_names):
1665     print(f"{class_name}: {cm_normalized[i][i]:.2%} correct")
1666
1667
1668 # #### F1, Precision, and Recall
1669
1670 # In[37]:
1671
1672
1673 from sklearn.metrics import classification_report
1674
1675 # Predict on the validation set using the best DNN model
1676 if best_model_state is not None:
1677     model.load_state_dict(best_model_state)
1678 model.eval()
1679 with torch.no_grad():
1680     X_val_tensor = torch.tensor(X_val_scaled, dtype=torch.float32).to(device)
1681     logits = model(X_val_tensor)
1682     y_val_pred = torch.argmax(logits, dim=1).cpu().numpy()
1683
1684 # Use the same label encoder as during training for class names
1685 class_names = le.classes_
1686
1687 # Print precision, recall, and F1 score for each class
1688 print("Precision, Recall, and F1 Score (Validation Set):\n")
1689 print(classification_report(y_val_enc, y_val_pred, target_names=class_names))

```