# USER DEPENDENT ANALYSIS

ASSIGNMENT #2

Professor Ayan Banarjee, Ira A. Fulton School of Engineering, ASU
Vanessa Ulloa | CSE 572: Data Mining | 28 Apr 2019

# TABLE OF CONTENTS

## PART 1: ASSIGNMENT 1

### PHASE 1:

In the first phase of the assignment we were given raw data for a select number of users that consisted of sensor data. There was IMU sensor data and EMG sensor data available for actions with either a fork or a spoon. Along with this data there was information provided regarding the video frames that corresponded with the separate eating and non-eating actions. This was used during the synchronization of the data in order to have ground truth data. The IMU data contained 10 different sensor information (OriX, OriY, OriZ, OriW, AccX, AccY, AccZ, GyroX, GyroY, GyroZ.) and the EMG was not utilized for this project due to poor synchronization. There were many issues in synchronizing the EMG data with the video frame data, this stems from the sampling rates for EMG data being inconsistent and therefore we cannot guarantee accuracy for EMG ground truth data.

### PHASE 2:

#### SYNCHRONIZATION

The synchronization of the IMU data required the assumptions that the frames per second was 30fps and that the sampling rate was 50 Hz or 50 samples per second. The data was synchronized using these assumptions and we were able to separate the Eating actions using the start and end frames and the Non-Eating actions (the actions in between the end frame and the next start frame).

This produced a Matrix where all the eating actions were saved. Each column in the matrix in the .csv file is representative of a different sensor mentioned before.

| Name | Status | Date modified |
|---|---|---|
| ☐ user10_IMU_Eat.csv | ⟳ | 4/7/2019 4:07 |
| user10_IMU_NotEat.csv | ⟳ | 4/7/2019 4:07 |
| user11_IMU_Eat.csv | ⟳ | 4/7/2019 4:07 |
| user11_IMU_NotEat.csv | ⟳ | 4/7/2019 4:07 |
| user12_IMU_Eat.csv | ⟳ | 4/7/2019 4:07 |
| user12_IMU_NotEat.csv | ⟳ | 4/7/2019 4:07 |
| user13_IMU_Eat.csv | ⟳ | 4/7/2019 4:07 |
| user13_IMU_NotEat.csv | ⟳ | 4/7/2019 4:07 |
| user14_IMU_Eat.csv | ⟳ | 4/7/2019 4:07 |
| user14_IMU_NotEat.csv | ⟳ | 4/7/2019 4:07 |
| user16_IMU_Eat.csv | ⟳ | 4/7/2019 4:07 |
| user16_IMU_NotEat.csv | ⟳ | 4/7/2019 4:07 |
| user17_IMU_Eat.csv | ⟳ | 4/7/2019 4:07 |
| user17_IMU_NotEat.csv | ⟳ | 4/7/2019 4:07 |

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | 0.779 | 0.565 | 0.083 | -0.259 | 0.895 | -0.3 | -0.371 | -1.312 |
| 2 | 0.779 | 0.565 | 0.086 | -0.26 | 0.84 | -0.261 | -0.354 | 15.25 |
| 3 | 0.777 | 0.567 | 0.09 | -0.258 | 0.841 | -0.375 | -0.354 | 33.688 |
| 4 | 0.776 | 0.57 | 0.095 | -0.254 | 0.87 | -0.504 | -0.3 | 34.25 |
| 5 | 0.773 | 0.574 | 0.098 | -0.251 | 0.913 | -0.454 | -0.279 | 31.875 |
| 6 | 0.771 | 0.579 | 0.101 | -0.247 | 0.845 | -0.396 | -0.299 | 39.812 |
| 7 | 0.769 | 0.583 | 0.105 | -0.24 | 0.868 | -0.458 | -0.228 | 44.688 |
| 8 | 0.768 | 0.586 | 0.109 | -0.236 | 0.89 | -0.515 | -0.235 | 33.688 |
| 9 | 0.766 | 0.588 | 0.111 | -0.235 | 0.811 | -0.491 | -0.228 | 13.188 |
| 10 | 0.765 | 0.59 | 0.11 | -0.233 | 0.826 | -0.483 | -0.209 | 1.75 |
| 11 | 0.766 | 0.592 | 0.105 | -0.23 | 0.818 | -0.479 | -0.132 | -11.188 |
| 12 | 0.768 | 0.591 | 0.098 | -0.226 | 0.882 | -0.418 | -0.136 | -25.688 |
| 13 | 0.77 | 0.589 | 0.09 | -0.226 | 0.854 | -0.307 | -0.213 | -39 |
| 14 | 0.774 | 0.585 | 0.083 | -0.227 | 0.805 | -0.243 | -0.289 | -37.5 |
| 15 | 0.778 | 0.582 | 0.076 | -0.225 | 0.889 | -0.245 | -0.355 | -26.438 |
| 16 | 0.781 | 0.579 | 0.07 | -0.221 | 0.884 | -0.385 | -0.257 | -23.062 |
| 17 | 0.785 | 0.576 | 0.065 | -0.22 | 0.935 | -0.477 | -0.159 | -33.125 |
| 18 | 0.787 | 0.573 | 0.058 | -0.221 | 0.959 | -0.34 | -0.243 | -45 |
| 19 | 0.789 | 0.569 | 0.051 | -0.224 | 0.866 | -0.236 | -0.336 | -43.625 |
| 20 | 0.792 | 0.566 | 0.045 | -0.225 | 0.833 | -0.221 | -0.329 | -32.938 |

In this diagram, the columns are related to the IMU data:

- Column A: Orientation X
- Column B: Orientation Y
- Column C: Orientation Z
- Column D: Orientation W
- Column E: Acceleration X
- Column F: Acceleration Y
- Column G: Acceleration Z
- Column H: Gyroscope X
- Column I: Gyroscope Y
- Column J: Gyroscope Z

## FEATURE EXTRACTION

The Feature exaction methods that were used were:

1. Mean
2. Standard Deviation
3. Range
4. Minimum

5. Maximum

The data obtained during feature selection resulted in a data set for all users used in the previous phase. PCA or Principal Component Analysis is a method of linear dimensional reduction. A new feature matrix was obtained for Eating and Non-Eating actions that can be used to train and test different classification models:

1. Decision Trees
2. Support Vector Machines
3. Neural Network Machines

The data was divided randomly (by eating and non-eating action) into 60% training data and 40% test data. Target data was produced by creating data matrices which the same dimensions as the test and target data and classifying with a 1 or 0 whether the target was an eating or non-eating action.

## USER DEPENDENT ANALYSIS

In the creation of the models for user analysis, the user data was compiled after PCA and the new feature matrix but still separated by Eating and Non-Eating actions. This data was randomly split up into 7 groups where each group had training data (60%) and test data (40%). The data was not individual split amongst each user due to small data sets not producing accurate results in the models and this gave us an unbiased view of data amongst many users but randomized to simulate an individual.
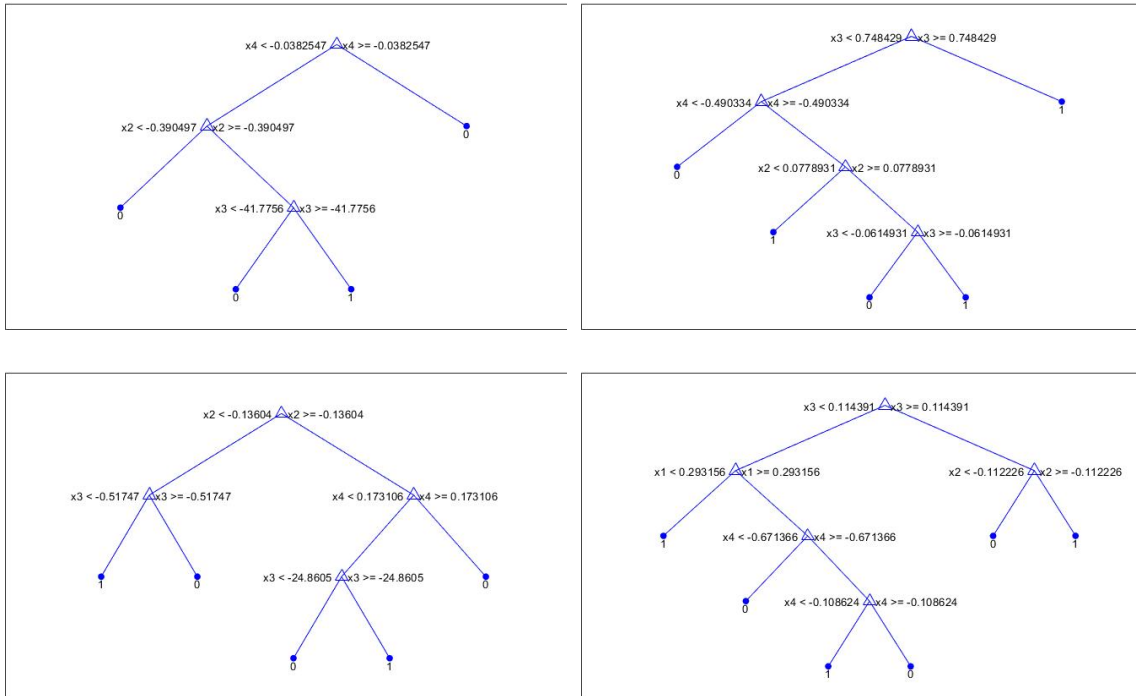
### DECISION TREE

The decision tree model was run against the seven different groups of data with test and train data for eating and non-eating actions along with target data that was used to train the model.
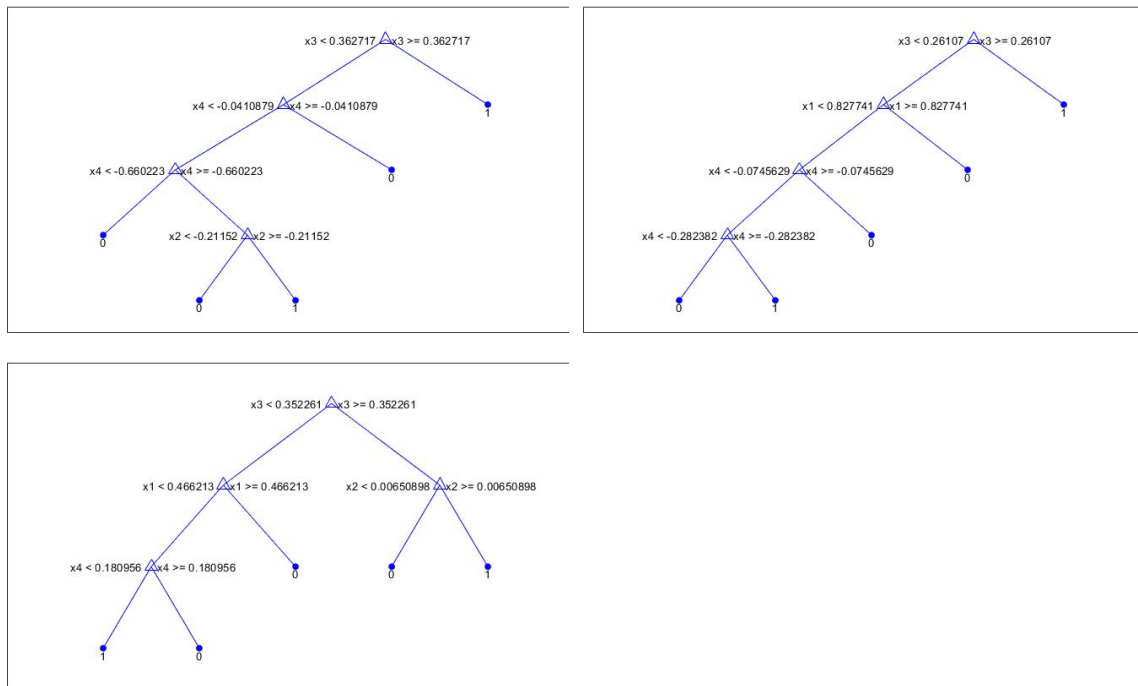
This produced 7 individual Precision, Recall and F1 Scores:

| Group | Decision Tree F1Score | Decision Tree Recall | Decision Tree Precision |
|--------|--------|--------|--------|
| Group1 | 0.5 | 0.97059 | 0.66 |
| Group2 | 0.34375 | 0.7915 | 0.47933 |
| Group3 | 0.5 | 0.94444 | 0.65385 |
| Group4 | 0.5 | 0.86364 | 0.63333 |
| Group5 | 0.28125 | 0.7 | 0.40127 |
| Group6 | 0.34375 | 0.75397 | 0.47221 |
| Group7 | 0.34375 | 0.83333 | 0.48673 |

The Precision score and Recall score varied for each group selected, the low F1 scores indicate that while for some groups the classification of positive observations was mostly correct, over all the misclassification percentage was high.

Decision Trees:

Tree 1:
x4 < -0.0382547 | x4 >= -0.0382547
x2 < -0.390497 | x2 >= -0.390497 → 0
0
x3 < -41.7756 | x3 >= -41.7756
0 | 1

Tree 2:
x3 < 0.748429 | x3 >= 0.748429 → 1
x4 < -0.490334 | x4 >= -0.490334
0
x2 < 0.0778931 | x2 >= 0.0778931
1
x3 < -0.0614931 | x3 >= -0.0614931
0 | 1

Tree 3:
x2 < -0.13604 | x2 >= -0.13604
x3 < -0.51747 | x3 >= -0.51747
1 | 0
x4 < 0.173106 | x4 >= 0.173106
x3 < -24.8605 | x3 >= -24.8605
0 | 1
0

Tree 4:
x3 < 0.114391 | x3 >= 0.114391
x1 < 0.293156 | x1 >= 0.293156
1
x4 < -0.671366 | x4 >= -0.671366
0
x4 < -0.108624 | x4 >= -0.108624
1 | 0
x2 < -0.112226 | x2 >= -0.112226
0 | 1

x3 < 0.362717 / x3 >= 0.362717
x4 < -0.0410879 / x4 >= -0.0410879
x4 < -0.660223 / x4 >= -0.660223
x2 < -0.21152 / x2 >= -0.21152
0   0   0   0   1

x3 < 0.26107 / x3 >= 0.26107
x1 < 0.827741 / x1 >= 0.827741
x4 < -0.0745629 / x4 >= -0.0745629
x4 < -0.282382 / x4 >= -0.282382
0   1   0   0   1

x3 < 0.352261 / x3 >= 0.352261
x1 < 0.466213 / x1 >= 0.466213
x2 < 0.00650898 / x2 >= 0.00650898
x4 < 0.180956 / x4 >= 0.180956
1   0   0   0   1

## SUPPORT VECTOR MACHINE

The support vector model was run against the seven different groups of data with test and train data for eating and non-eating actions along with target data that was used to train the model.

| Group | SVM F1Score | SVM Recall | SVM Precision |
|---|---|---|---|
| Group1 | 0.21875 | 0.75 | 0.33871 |
| Group2 | 0.46875 | 0.75 | 0.57692 |
| Group3 | 0.5 | 0.7963 | 0.61429 |
| Group4 | 0.3125 | 0.72672 | 0.43706 |
| Group5 | 0.25 | 0.83333 | 0.38462 |
| Group6 | 0.3125 | 0.72672 | 0.43706 |
| Group7 | 0.375 | 0.66194 | 0.47877 |

The F1 scores indicate there was poor precision and poor recall, overall for the model the precision or classification of true observations was relatively low when for all the groups the recall was relatively high. Due to F1 being a weighed average the low precision score would cause a low F1 score.

## NEURAL NET MACHINE

The Neural Net model (using *patternnet*) was run against the seven different groups of data with test and train data for eating and non-eating actions along with target data that was used to train

the model. During an observation made during the independent analysis the Percent Error of misclassification could be reduced by increasing the number of hidden neurons, so the increase was made from 10 hidden neurons to 30 for each group of data:

| Group | NN F1Score | NN Recall | NN Precision |
| --- | --- | --- | --- |
| Group1 | 0.61765 | 0.5 | 0.80769 |
| Group2 | 0.61429 | 0.5 | 0.7963 |
| Group3 | 0.60584 | 0.46875 | 0.85628 |
| Group4 | 0.60526 | 0.5 | 0.76667 |
| Group5 | 0.38462 | 0.25 | 0.83333 |
| Group6 | 0.60256 | 0.5 | 0.75806 |
| Group7 | 0.40796 | 0.28125 | 0.74242 |

The F1 score for each group was not very high and over all the Recall score was not very high either, the Precision score was consistently high however, but this still shows a high misclassification percentage.

Increasing the number of hidden neurons for the model did not change the results significantly:

| Group | NN F1Score | NN Recall | NN Precision |
| --- | --- | --- | --- |
| Group1 | 0.5 | 0.375 | 0.75 |
| Group2 | 0.19342 | 0.125 | 0.42727 |
| Group3 | 0.625 | 0.5 | 0.83333 |
| Group4 | 0.60526 | 0.5 | 0.76667 |
| Group5 | 0.60811 | 0.5 | 0.77586 |
| Group6 | 0.45132 | 0.34375 | 0.65686 |
| Group7 | 0.53645 | 0.4375 | 0.69324 |

## USER INDEPENDENT ANALYSIS

### DECISION TREE

The training data was passed to the decision tree model along with the training target data. This model received the test data and the result was a confusion matrix.

| 1 | 2 |
| --- | --- |
| 102 | 9 |
| 10 | 103 |

The confusion matrix can be used to calculate the precision and recall, which can in turn be used to calculate the f-score.

| Decision Tree Results | |
|---|---|
| Precision | 0.9152 |
| Recall | 0.4558 |
| F1-Score | 0.6085 |

Precision is the calculation of total True Positive divided by the Sum of True Positive and False Positive. This calculation measures how accurate this decision tree model is when predicting positive results. In this decision tree model, the score was 0.9152.

Recall is the calculation of true positives being calculated through the model. In this decision tree model, the score was 0.4558.

F1 score is used when a balance is required between precision and recall, these are measures relating to the number of observations that were correctly classified as positive versus the expected positive observations. The F1 score is used because it considers false positives and false negatives (or incorrect classifications). The F1 score for this model is 0.6085.

The decision tree produced by this model is below



## SUPPORT VECTOR MACHINE

The same test and training data were passed to the SVM or Support Vector Machine model, this was the resulting confusion matrix:

| | |
|---|---|
| 98 | 44 |
| 14 | 68 |

In this format we can use the confusion matrix to calculate the Precision, Recall and F1 Scores for this classification model:

| SVM Results | |
|---|---|
| Precision | 0.7597 |
| Recall | 0.3036 |
| F1-Score | 0.4338 |

In this model the 0.7597 score for Precision is a pretty good score. Precision is the measure of how many positive observations were correctly classified against false positives. The Recall score was a bit low at 0.3036 which means that this measure of positive observations correctly classified against that were labeled correctly is low. Due to the lower Recall score the F-Score also came in low since this score is the weighted average and considers false negatives and true negatives.

## NEURAL NET MACHINE

The same Training and Test data was passed to the a *patternnet* (Neural Net used for classification) with 10 hidden neurons, which resulted:

Neural Pattern Recognition (nprtool)    —    □    ✕

## Train Network

Train the network to classify the inputs according to the targets.

### Train Network

Train using scaled conjugate gradient backpropagation.    (trainscg)

[ 🐦 Retrain ]

Training automatically stops when generalization stops improving, as indicated by an increase in the cross-entropy error of the validation samples.

### Results

| | Samples | CE | %E |
|---|---|---|---|
| Training: | 201 | 7.76850e-1 | 36.31840e-0 |
| Validation: | 17 | 1.03031e-0 | 23.52941e-0 |
| Testing: | 118 | 8.51044e-1 | 39.83050e-0 |

[ Plot Confusion ]    [ Plot ROC ]

### Notes

🔸 Training multiple times will generate different results due to different initial conditions and sampling.
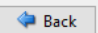
📧 Minimizing Cross-Entropy results in good classification. Lower values are better. Zero means no error.
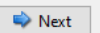
% Percent Error indicates the fraction of samples which are misclassified. A value of 0 means no misclassifications, 100 indicates maximum misclassifications.
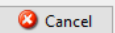
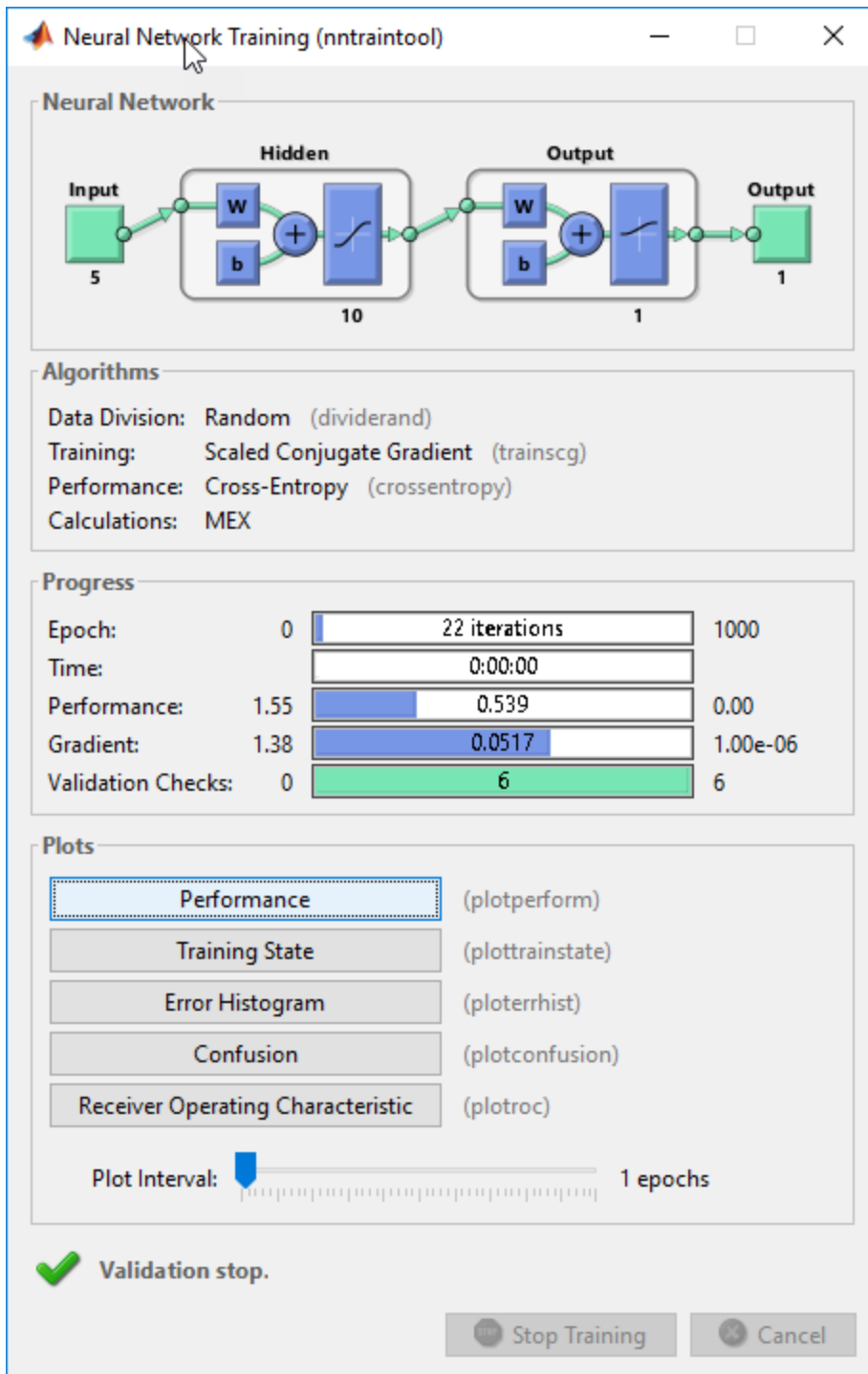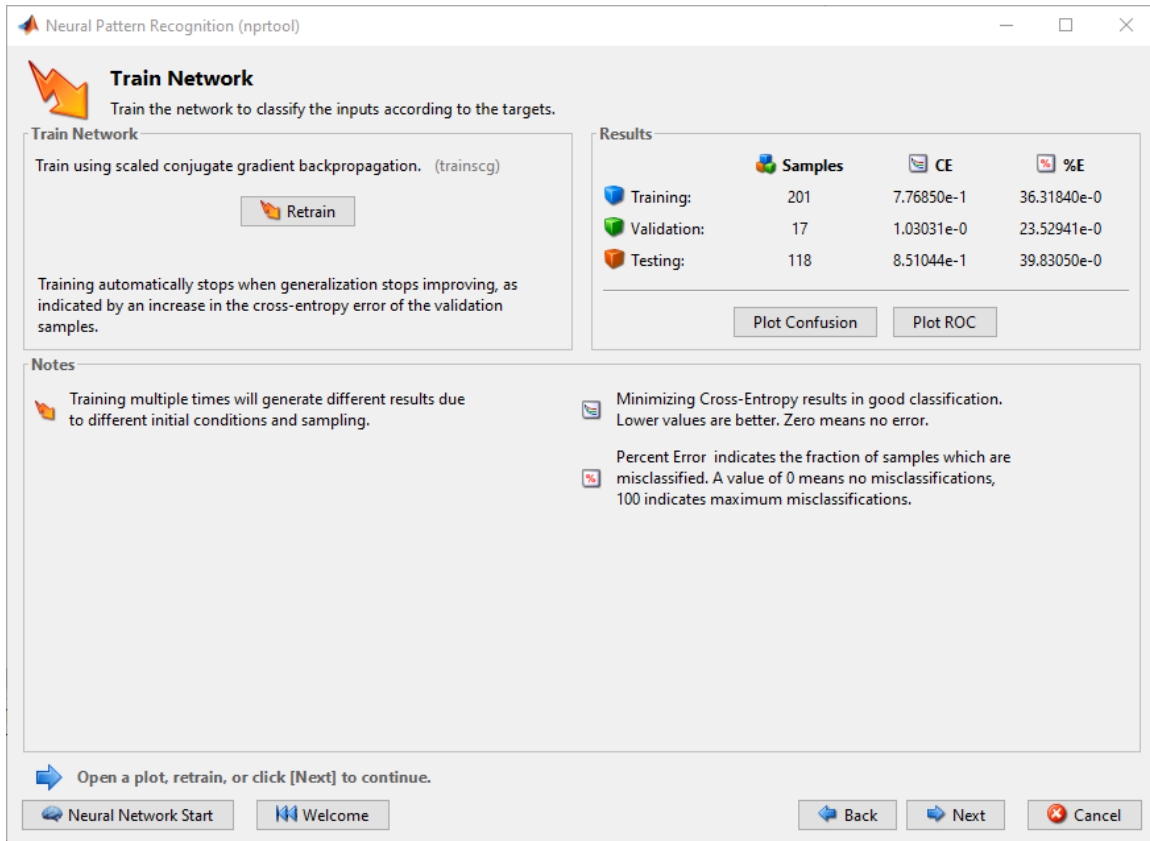➡ Open a plot, retrain, or click [Next] to continue.

[ 🌐 Neural Network Start ]    [ ◀◀ Welcome ]    [ ⬅ Back ]    [ ➡ Next ]    [ ✖ Cancel ]

The model had a high level of Percent Error which represents how many samples were misclassified.

When the number of hidden neurons was increase from 10 to 39 neurons, the Percent Error decreased from 36% during Training to 27% and in the Test Data from 39% to 24%.

Using the new neural network configuration of 30 hidden neurons, the test and target test data was passed through the model and resulted in 24% Percent Error (misclassification).

| NN Results | |
|---|---|
| Precision | 0.7250 |
| Recall | 0.3839 |
| F1-Score | 0.5020 |

The low F1 score matches the observations made earlier with the high classification errors, and it is also seen the low Recall score. The precision score was high which indicates that the classification of positive observations was less effected.

## CONCLUSION

In the Dependent and Independent analysis of this data using Support Vector Machines, Decision Trees and Neural Network Machines the Precision, Recall and F1-scores varied depending on the amount of test and training data. Our observation is that the model data was affected by the feature selection and PCA results and since the EMG data was removed due to its inconsistency. The variance that the EMG data would have provided could have shown better results for the binary classification models.