

CSE 572 Data Mining Assignment – 4

Team bit_miner's

Ankit Anand 1213250712

Rohith Reddy Vajrala 1212908292

Lakshmisagar Kusnoor 1211009498

Divya Prakash Sivakumar 1213204601

We used the PCA and features set from previous assignment as new feature set by multiplying both and divided this new feature set into two parts for each user:

Part 1: Training and

Part 2: Test.

We used 60% of the data for each user as training and the rest of 40% as test data.

We used three types of machines in this assignment:

- a) **Decision trees,**
- b) **Support vector machines and**
- c) **Neural networks.**

We trained each machine with the training data and then used the test data to report accuracy.

Using the accuracy metrics of Precision, Recall, F1 score and ROC reported each metric for every group.

Following are the Formulae's Used:

1. **Accuracy:** is the ratio of correct predictions to total predictions made.

$$\text{classification accuracy} = \text{correct predictions} / \text{total predictions} * 100$$

2. **Accuracy matrix/ Confusion Matrix:** A confusion matrix is a summary of prediction results on a classification problem. It shows the ways in which our classification model is confused when it makes predictions.

Process for calculating an Accuracy Matrix

- We used our test data with expected outcome values.
- Made predictions for each row in our test dataset.
- From the expected outcomes and predictions count:
 - The number of correct predictions for each class.
 - The number of incorrect predictions for each class, organized by the class that was predicted.
- These numbers are then organized into a matrix with each row of the matrix corresponds to a predicted class and each column of the matrix corresponds to an actual class.
- The counts of correct and incorrect classification are then filled into the matrix.

Actual Class	Predicted class		
		Class = Yes	Class = No
	Class = Yes	True Positive	False Negative
	Class = No	False Positive	True Negative

3. **Precision:** It tells what percent of positive predictions were correct. It is the ratio of correctly predicted positive observations to the total predicted positive observations

$$\text{Precision} = \text{True Positive} / \text{True Positive} + \text{False Positive}$$

4. **Recall:** Tells what percent of the positive cases did we catch. It is the ratio of correctly predicted positive observations to the all observations in actual class

$$\text{Recall} = \text{True Positive} / \text{True Positive} + \text{False negative}$$

5. **F1 Score:** F1 Score is the weighted average of Precision and Recall.

$$\text{F1 Score} = 2 * (\text{Recall} * \text{Precision}) / (\text{Recall} + \text{Precision})$$

6. **ROC Curve:** A Receiver Operating Characteristic (ROC) Curve is a way to compare diagnostic tests. It is a plot of the true positive rate against the false positive rate.

7. **AUC:** Area under the curve(AUC) is equal to the probability that a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one (assuming 'positive' ranks higher than 'negative').

1. **Decision tree (Classification):** A decision tree is a decision support tool that uses a tree-like graph or model of decisions and their possible consequences, including chance event outcomes, resource costs, and utility. It is one way to display an algorithm that only contains conditional control statements.

Procedure:

- We read the Eating and Non- Eating csv files
- Then we separate the above data as training and testing
- We pass the train data as parameter to **fitctree** method of MATLAB which returns a classifier. This classifier is further used to predict the scores against the test data.
- Now that we have eating and non-eating predictions. We calculate True Positive matrix, False Negative matrix, True Negative matrix, False Positive matrix.
- Using the above matrices, we calculate Precision, Recall, F1 Score and AUC and write the result into a file.

Classification decision tree gives responses in the leaf node in form of true or false.

Dimensions of feature matrix:

Training Data: 90 x 1440 (60% of this)

Testing Data: 90 x 1440 (40% of this)

Class-Name (Labels): 1 (Eating actions) and 0 (Non-Eating action)

2. **Support vector machine:** A Support Vector Machine (SVM) is a discriminative classifier formally defined by a separating hyperplane. In other words, given labeled training data (supervised learning), the algorithm outputs an optimal hyperplane which categorizes new examples.

Procedure:

- We read the Eating and Non- Eating csv files
- Then we separate the above data as training and testing
- We pass the train data as parameter to **fitcsvm** method of MATLAB which returns a support vector machine classifier. This classifier is further used to predict the scores against the test data.
- Now that we have eating and non-eating predictions. We calculate True Positive matrix, False Negative matrix, True Negative matrix, False Positive matrix.
- Using the above matrices, we calculate Precision, Recall, F1 Score and AUC and write the a result file.

Dimensions of feature matrix:

Training Data: 90 x 1440 (60% of this)

Testing Data: 90 x 1440 (40% of this)

Class-Name (Labels): 1 (Eating actions) and 0 (Non-Eating action)

3. **Neural Network:** is a learning algorithm that is inspired by the structure and functional aspects of biological neural networks. They are usually used to model complex relationships between inputs and outputs, to find patterns in data, or to capture the statistical structure in an unknown joint probability distribution between observed variables.

Procedure:

- We read the Eating and Non- Eating csv files
- Then we separate the above data with training and testing
- We pass the train data as parameter to **fitnet** method which has parameter of 10 neurons and trainlm algorithm (trainlm is a network training function that updates weight and bias values according to Levenberg-Marquardt optimization.). Then we train the machine and predict the scores against the test data.
- Now that we have eating and non-eating predictions. We calculate True Positive matrix, False Negative matrix, True Negative matrix, False Positive matrix.
- Using the above matrices, we calculate Precision, Recall, F1 Score and AUC and write a result file.

Dimensions of feature matrix:

Training Data: 90 x 1440 (60% of this)

Testing Data: 90 x 1440 (40% of this)

Class-Name (Labels): 1 (Eating actions) and 0 (Non-Eating action)

=====

Phase 1: User dependent analysis result

We have combined each group’s spoon and fork data for both the categories and have randomly selected 60% (random sampling) of the data as train data and the rest was used as test data to report the accuracy metrics such as Precision, Recall, F1 score and ROC for each of the three models, which are Decision Trees (fitctree) and SVM (fitsvm) and Neural Network (Neural Network Toolbox).

Group	Decision Tree				Neural Network				SVM			
	Precision	Recall	F1	AUC	Precision	Recall	F1	AUC	Precision	Recall	F1	AUC
1	1	0.965517	0.982456	0.985714	0.925926	0.862069	0.892857	0.889286	1	0.967742	0.983607	0.984848
2	0.941176	0.969697	0.955224	0.928922	0.878788	0.878788	0.878788	0.796537	1	1	1	1
3	0.931034	0.9	0.915254	0.9214	0.777778	0.7	0.736842	0.763889	1	1	1	1
4	0.962963	0.866667	0.912281	0.925926	0.83871	0.866667	0.852459	0.841734	1	0.972222	0.985915	0.982143
5	0.857143	1	0.923077	0.928571	0.727273	0.8	0.761905	0.763636	1	1	1	1
6	1	0.969697	0.984615	0.983871	0.892857	0.757576	0.819672	0.832143	1	1	1	1
7	0.884615	0.958333	0.92	0.928794	0.676471	0.958333	0.793103	0.810714	0.962963	1	0.981132	0.981481
8	0.96875	0.96875	0.96875	0.968246	0.96875	0.96875	0.96875	0.968246	1	1	1	1
9	0.777778	1	0.875	0.888889	0.777778	1	0.875	0.888889	1	1	1	1
10	0.933333	0.777778	0.848485	0.84902	0.666667	0.533333	0.592593	0.658333	1	1	1	1
11	1	0.903226	0.949153	0.957143	0.774194	0.774194	0.774194	0.777722	1	1	1	1
12	0.742857	0.866667	0.8	0.784472	0.675	0.870968	0.760563	0.732237	0.928571	0.896552	0.912281	0.914286
13	1	1	1	1	0.333333	0.266667	0.296296	0.391667	1	1	1	1
14	0.807692	0.7	0.75	0.782225	0.75	0.617647	0.677419	0.656187	1	0.914286	0.955224	0.951613
15	0.820513	1	0.901408	0.910256	0.764706	0.8125	0.787879	0.778905	0.961538	0.925926	0.943396	0.953742
16	0.722222	0.928571	0.8125	0.824074	0.714286	0.535714	0.612245	0.702381	0.928571	0.83871	0.881356	0.892857
17	0.925926	0.78125	0.847458	0.865741	0.848485	0.933333	0.888889	0.923935	1	0.965517	0.982456	0.985714
18	0.909091	1	0.952381	0.954545	0.807692	0.7	0.75	0.782225	0.941176	1	0.969697	0.970588
19	0.866667	0.962963	0.912281	0.918182	0.807692	0.777778	0.792453	0.822765	0.966667	1	0.983051	0.983333
20	0.933333	0.933333	0.933333	0.904167	0.742857	0.866667	0.8	0.68961	0.967742	1	0.983607	0.983871
21	0.96875	1	0.984127	0.984375	0.88	0.6875	0.77193	0.808421	1	1	1	1
22	0.787234	0.948718	0.860465	0.831117	0.961538	0.641026	0.769231	0.773774	0.967742	1	0.983607	0.983871

23	0.93333 3	0.9032 26	0.9180 33	0.9212 12	1	0.8064 52	0.8928 57	0.9210 53	0.93333 3	1	0.9655 17	0.9666 67
24	1	0.9375	0.9677 42	0.9696 97	0.77419 4	0.75	0.7619 05	0.7620 97	1	0.9714 29	0.9855 07	0.9827 59
25	1	0.9705 88	0.9850 75	0.9827 59	0.83783 8	0.8857 14	0.8611 11	0.8419 96	0.90909 1	1	0.9523 81	0.9545 45
26	0.89189 2	0.9705 88	0.9295 77	0.9267 15	0.74074 1	0.5882 35	0.6557 38	0.6759 26	0.94117 6	0.9142 86	0.9275 36	0.9188 64
27	0.83871	0.9629 63	0.8965 52	0.8966 28	0.95	0.6785 71	0.7916 67	0.8426 47	0.95	0.95	0.95	0.9598 48
28	0.95454 5	0.75	0.84	0.8919 07	0.85714 3	0.75	0.8	0.8142 86	1	0.9375	0.9677 42	0.9696 97
29	1	1	1	1	1	0.9117 65	0.9538 46	0.9531 25	1	1	1	1
30	0.96969 7	1	0.9846 15	0.9848 48	0.86363 6	0.5757 58	0.6909 09	0.6136 36	0.92857 1	0.9629 63	0.9454 55	0.9309 52
31	0.64	0.7619 05	0.6956 52	0.5422 22	0.72222 2	0.5416 67	0.6190 48	0.5173 61	0.95	0.7916 67	0.8636 36	0.7964 29
32	1	1	1	1	0.55319 1	1	0.7123 29	0.7765 96	1	1	1	1
33	0.84375	1	0.9152 54	0.9218 75	0.78571 4	0.6111 11	0.6875	0.6734 28	0.97297 3	1	0.9863 01	0.9864 86
Average	0.9034 24	0.9290 28	0.9127 5	0.9110 16	0.7962 26	0.7548 12	0.7660 6	0.7710 72	0.9760 64	0.9699 64	0.9724 06	0.9707 45

Analysis of Phase 1 Result: - We think that in phase 1 since we are training every group with their own new feature set and as the training data is 60% of the feature set, thus the machine is getting overfit and giving very high accuracy and all the other class metrics are also high.

Since the machine is trained well with good percentage of training data from every group's dataset we are getting high classification accuracy results.

For this phase Decision tree and Support Vector machine performed better than Neural Network.

Phase 2: User independent analysis result

We combined the first 10 groups data as the train data both the categories of spoon and fork. The data of the other 23 groups were selected without combining spoon and fork data and used to test the models and report the accuracy metrics

Group	Decision Tree				Neural Network				SVM			
	Precision	Recall	F1	AUC	Precision	Recall	F1	AUC	Precision	Recall	F1	AUC
11	0.65	0.4632 5	0.4698	0.5961 86	0.53488 4	0.575	0.5542 17	0.5313 31	0.50236	0.6123	0.5936	0.5311
12	0.66666 7	0.3297 3	0.4269 23	0.6397 44	0.45945 9	0.2297 3	0.3063 06	0.4706 39	0.44898	0.4556 96	0.4864 86	0.4705 88
13	0.68627 5	0.875	0.7692 31	0.7569 3	0.57575 8	0.9743 59	0.7238 1	0.7066 59	0.42663	0.3945 6	0.4102 5	0.4126 3

14	0.65789 5	0.3125	0.4237 29	0.5997 81	0.57142 9	0.55	0.5605 1	0.5634 92	0.50632 9	0.5111 11	0.575	0.5411 76
15	0.61	0.5569	0.4333 33	0.5957 63	0.57471 3	0.625	0.5988 02	0.5760 89	0.57594 9	0.5822 78	0.575	0.5786 16
16	0.64864 9	0.523	0.4102 56	0.5929 19	0.48	0.6	0.5333 33	0.4641 38	0.46202 5	0.4712 64	0.5125	0.4910 18
17	0.54902	0.469	0.4274 81	0.5315 19	0.59154 9	0.525	0.5562 91	0.5773 84	0.54256	0.5136 4	0.5010 6	0.4956 8
18	0.56488	0.556	0.514	0.4989 8	0.63414 6	0.325	0.4297 52	0.5863 04	0.56962	0.5652 17	0.65	0.6046 51
19	0.5625	0.4623	0.4923 65	0.5312 5	0.54386	0.3875	0.4525 55	0.5293 56	0.44303 8	0.4534 88	0.4875	0.4698 8
20	0.78378 4	0.3670 89	0.6512 46	0.5793 92	0.83333 3	0.5063 29	0.6299 21	0.6340 58	0.47008 5	0.6349 21	0.5063 29	0.5633 8
21	0.621	0.3956 4	0.4638 5	0.498	0.51219 5	0.5316 46	0.5217 39	0.5094 31	0.45222 9	0.4578 31	0.4810 13	0.4691 36
22	0.79487 2	0.3827 16	0.5166 67	0.6891 03	0.63953 5	0.6790 12	0.6586 83	0.6416 85	0.62893 1	0.6341 46	0.6419 75	0.6380 37
23	0.75	0.6759 48	0.6123 5	0.6631 36	0.60714 3	0.85	0.7083 33	0.6731 37	0.50632 9	0.5116 28	0.55	0.5301 2
24	0.67857 1	0.5924 3	0.4523 6	0.6046 7	0.55	0.55	0.55	0.5442 31	0.4956	0.4321	0.4625	0.4836
25	0.65306 1	0.5236	0.4961 24	0.6043 08	0.66666 7	0.55	0.6027 4	0.6355 31	0.49044 6	0.5	0.475	0.4871 79
26	0.75264 1	0.6548 9	0.5946 1	0.6574	0.45679	0.4625	0.4596 27	0.4491 74	0.64525	0.6123 5	0.5823	0.5289 1
27	0.6965	0.6112	0.6523	0.6425	0.41666 7	0.3676 47	0.3906 25	0.4177 93	0.43283 6	0.4428 57	0.4558 82	0.4492 75
28	0.67	0.712	0.6945	0.6813	0.55172 4	0.8	0.6530 61	0.5853 86	0.38607 6	0.4086 02	0.475	0.4393 06
29	0.81818 2	0.6751 24	0.3529 41	0.6811 5	0.54929 6	0.975	0.7027 03	0.7121 48	0.71256	0.6123 5	0.6547	0.6248
30	0.75	0.657	0.4485 98	0.5438 31	0.69135 8	0.7466 67	0.7179 49	0.5063 93	0.49541 3	0.6612 9	0.5466 67	0.5985 4
31	0.85	0.7103	0.4415 58	0.6173 08	0.82352 9	0.7368 42	0.7777 78	0.6911 76	0.63529 4	0.7241 38	0.7368 42	0.7304 35
32	0.63265 3	0.5846 5	0.5564 8	0.5915 56	0.66666 7	0.625	0.6451 61	0.6526 1	0.43038	0.4390 24	0.45	0.4444 44
33	0.6423	0.7512	0.6812	0.756	0.66666 7	0.55	0.6027 4	0.6376 81	0.59631	0.5236 4	0.5452 3	0.5136

Average	0.68215	0.558325	0.520952	0.615336	0.59119	0.596619	0.579854	0.578079	0.515445	0.528454	0.537167	0.525918
---------	---------	----------	----------	----------	---------	----------	----------	----------	----------	----------	----------	----------

Analysis of Phase 1 Result: - We think that in phase 2 since we are training the machine with random 10 groups new feature set and testing on each group individually, the accuracy is low, and all the other class metrics are also low. This is because the machine will never be able to define a generalized classification function since the dataset for each group is different even though the feature vectors used are same, this is because each group eating and non-eating “actions/habits” are different and since the training has been done on different groups, it is tough to find a generalized machine.

Since the machine is trained with random 10 groups training data and being tested on each of the different groups individually, we are getting low classification accuracy results.

For this phase Decision tree performed better than Neural Network and Support Vector. Overall, we had much better classification accuracy results for Phase 1 then Phase 2.

=====

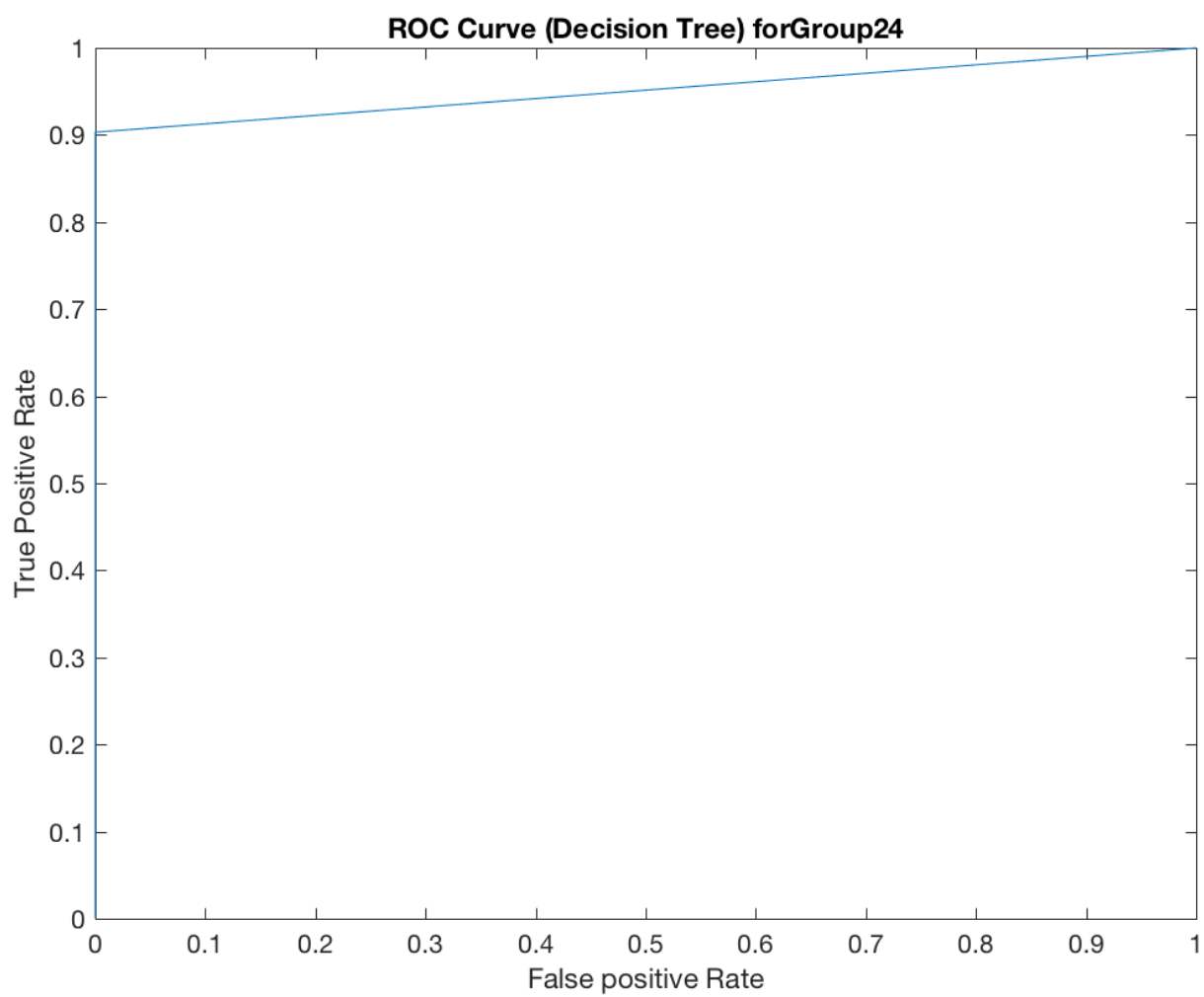
Sample ROC curves:

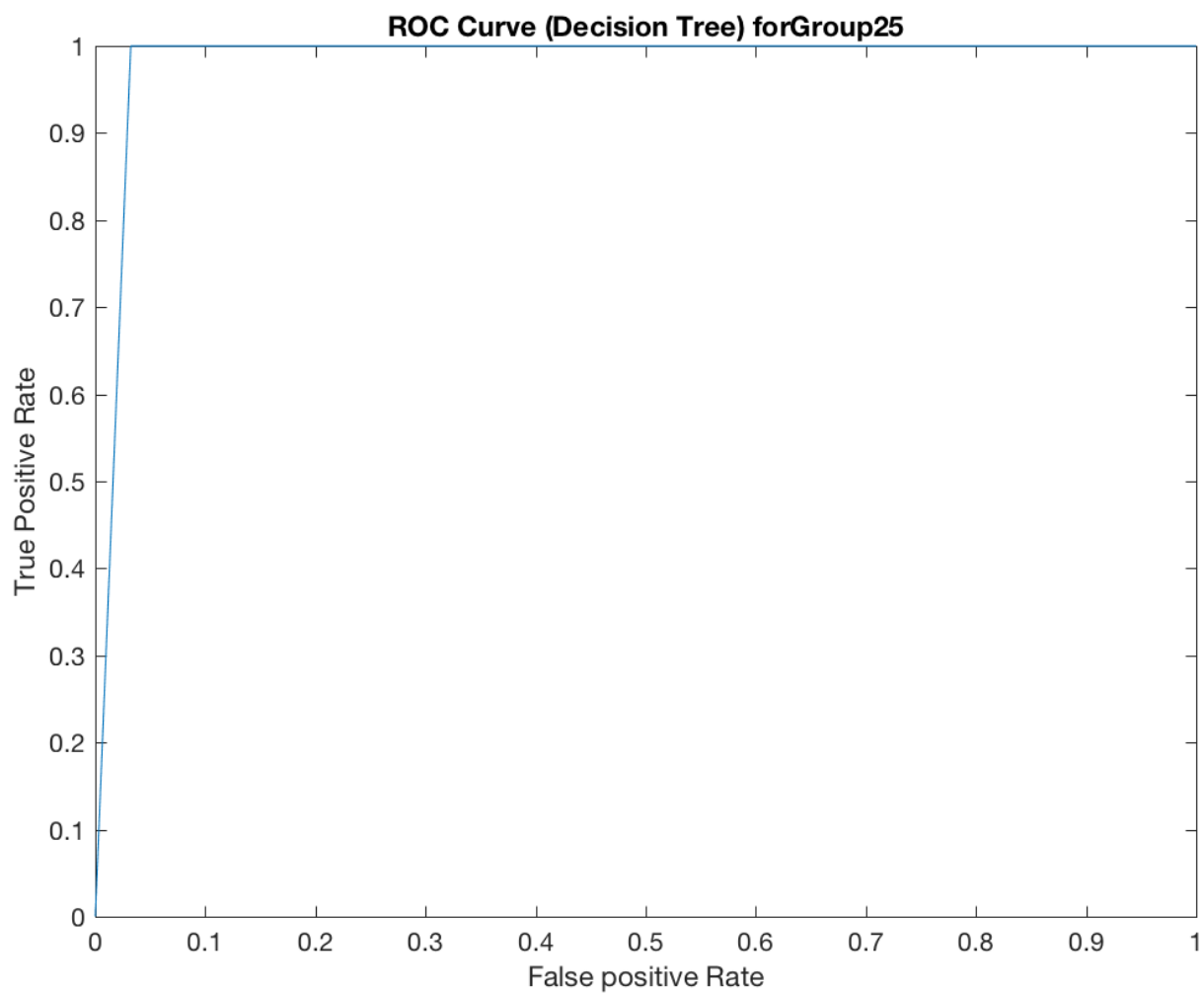
We have submitted ROC curves for 33 groups using each machine i.e. total 99 ROC curves in the submission zip file inside “ROC_Curves_**Phase 1**” folder. Their names follow the following pattern: -

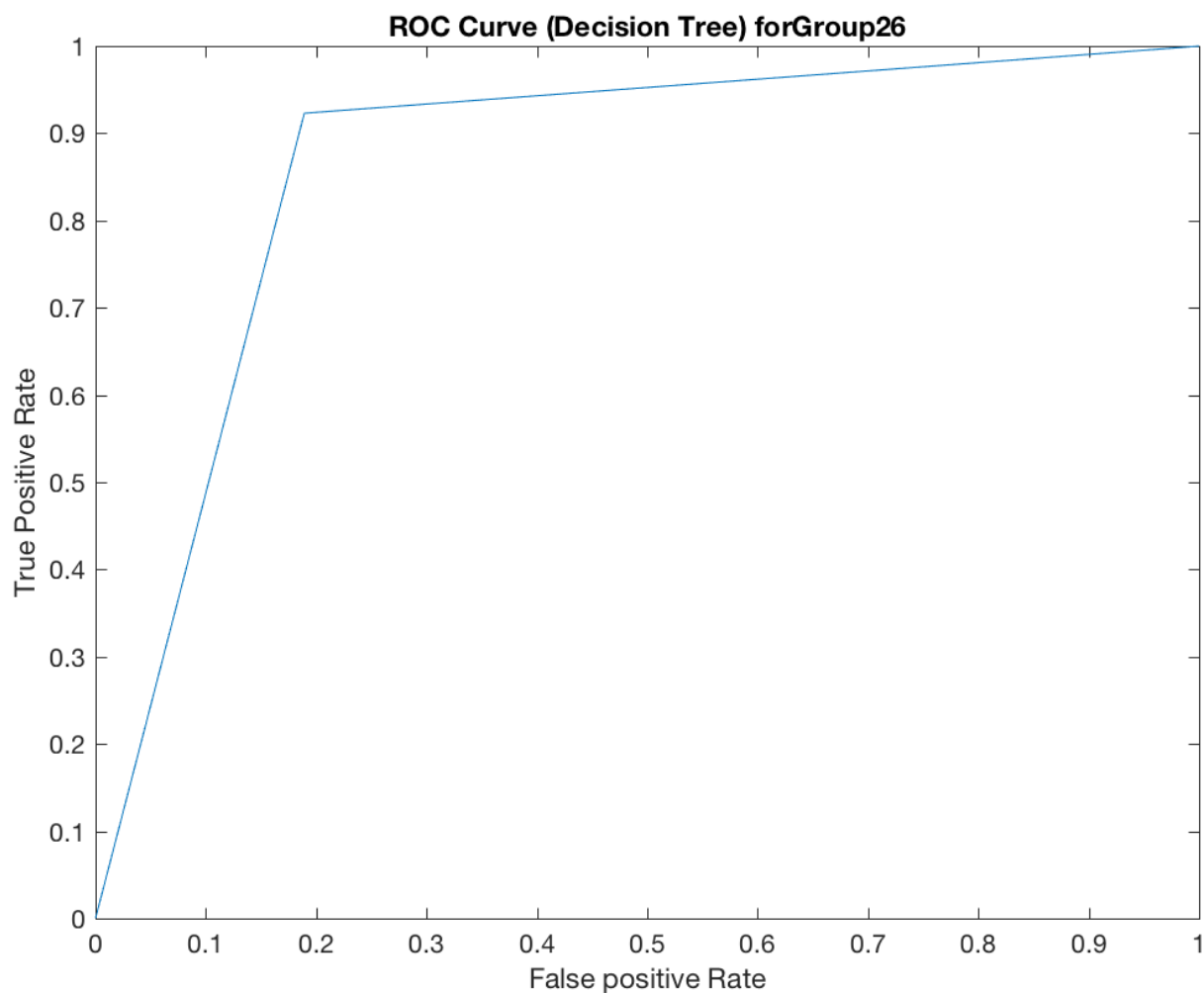
“RocCurve ”MachineType” GroupNo.”

We have submitted ROC curves for 23 groups using each machine i.e. total 69 ROC curves in the submission zip file inside “ROC_Curves_**Phase 2**” folder. Their names follow the following pattern: -

“RocCurve ”MachineType task2” GroupNo.”







Execution Steps(Readme):

1. Download the submitted Zip file and extract.
2. Find the final_task1.m file and run the file in MATLAB.
3. After execution completion, you will find the result.mat file which has 33 rows * 12 columns, which has Precision, Recall, F1 and AUC for three model decision tree , SVM and neural network ordered on the basis of group numbers.
4. Find the final_task2.m file for executing task 2 and run this file in MATLAB.
5. After execution you get resultmat_task2.mat file which has 23 rows* 12 columns of result.
6. Roc curves will be saved in the same folder path.

References:

- Wiki: https://en.wikipedia.org/wiki/Decision_tree 2.
- https://en.wikipedia.org/wiki/Artificial_neural_network
- <https://www.mathworks.com/help/stats/fitctree.html>
- <https://www.mathworks.com/help/stats/fitcsvm.html>
- https://en.wikipedia.org/wiki/Support_vector_machine