

# Deep Fakes Detection

Akul S. Malhotra<sup>1</sup> and Sana Naz<sup>1</sup>

<sup>1</sup>American University, College of Arts of Sciences, Washington, D.C., USA

## Introduction

Deep fakes leverage powerful techniques from machine learning and AI to manipulate or generate visual and audio content with a high potential to deceive. Large amount of readily accessible data and advanced computational power has made it easier to generate Deep fakes.

The motivation for this project is to build and evaluate an algorithm that detects if a particular media is fake or real. Accurate and efficient detection technology will prevent:

- Hampering societal peace
- Attacks on national democratic structure
- Prevent national and global security threat
- Stop the spread of fake news and its effects

## Contribution

We extract features from and build soft biometric models of individual's faces within a dataset of original and manipulated videos. Then, we evaluate the deep fake detection models in their ability to distinguish between real and fake videos and show the efficacy of this approach on a large amount of data.

## Classification Models

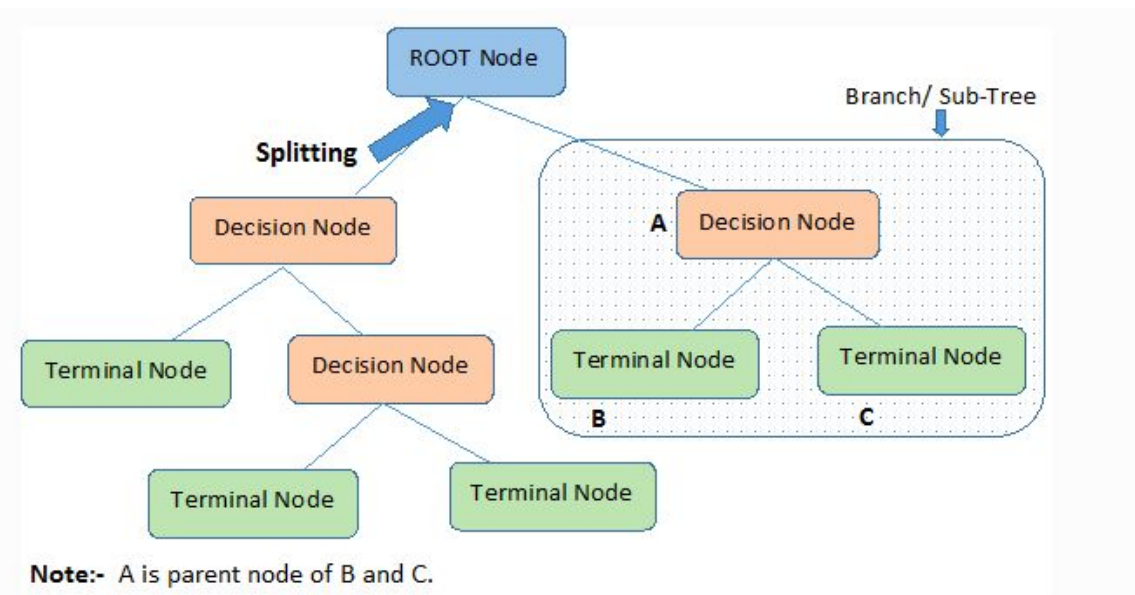
### Logistic Regression

This is used to predict the binomial outcome of a response variable using one or several predictor variables. The predictors can be binomial, categorical, or numerical. It is a way to map a continuous function of predictors to the probability of a binary response from 0 to 1.

$$\ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x$$
$$\Rightarrow P = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

### Decision Trees

The decision tree for a classification task is built through binary recursive partitioning- starting at the tree root and splitting the data into partitions on the feature that results in the largest information gain. Then, splitting it further on each of the nodes until the leaves are pure.



### Naive Bayes

This classifier is based on the Bayes theorem of probability to predict the class of the test dataset. It makes an assumption of independence amongst the predictors.

Bayes theorem provides a way of calculating posterior probability  $P(c/x)$  from  $P(c)$ ,  $P(x)$ , and  $P(x/c)$ :

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

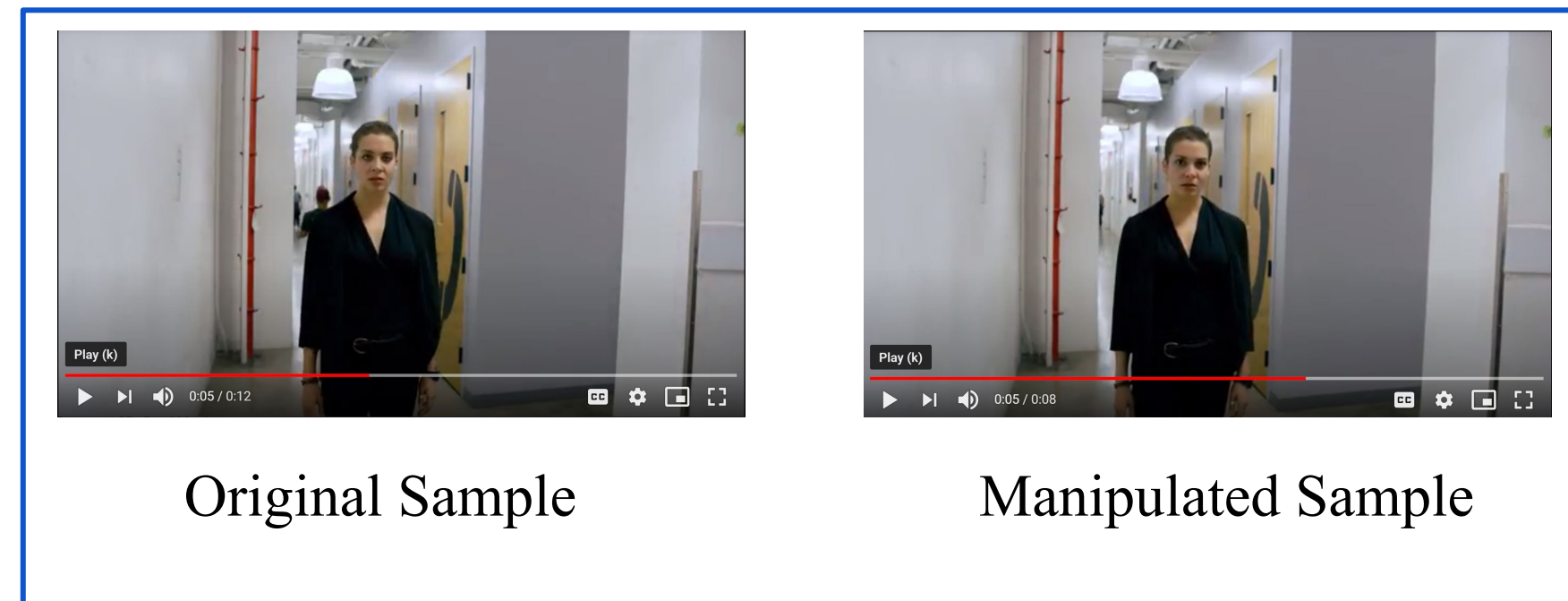
Labels: Likelihood, Class Prior Probability, Posterior Probability, Predictor Prior Probability

$$P(c|x) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

## Data Extraction Using OpenFace2

OpenFace2 is a forensic technique that is designed to detect deep fakes of an individual based on the principle that that when they speak, they exhibit relatively distinct patterns of facial and head movements that are disrupted by all 3 types of Deep fakes.

OpenFace2 is an open-source facial behavior analysis toolkit. All shortlisted videos contain 1 primary and camera-facing person in the frame.



We fed a sample of 150 original and manipulated videos to OpenFace2 and extract upper-body movements for speaker identification.

The output for each video are several features representing movements in each frame. Measurements like the following are later transformed into features:

- 2D and 3D facial landmark positions
- Head pose
- Eye gaze
- Facial action units (16)

## Data Organization & Feature Description

Each observation matrix  $X^{[k]}$  contains feature vectors which correspond to the movements of facial muscles commonly termed as Action Units (AU). In other words, each video contains frames at which certain facial expression were captured at a particular timestamp.

A total of 714 features were captured, for each video, using feature extraction methodology provided by OpenFace2.

Below is the hyperlink to follow for detailed information on the features metadata information:

<https://github.com/TadasBaltrusaitis/OpenFace/wiki/Output-Format>

Initial data exploration led us to incorporate a concise version of this dataset having 19 features.

Below is a snippet of cleaned data:

	AU01_r	AU02_r	AU04_r	AU05_r	AU06_r	AU07_r	AU09_r	AU10_r	AU12_r	AU14_r	AU15_r	AU17_r	AU2k
0	0.12	0.0	0.63	0.34	0.0	0.00	0.0	0.0	0.0	0.0	0.81	0.90	
1	0.12	0.0	0.78	0.39	0.0	0.05	0.0	0.0	0.0	0.0	0.76	0.77	
2	0.17	0.0	0.81	0.37	0.0	0.05	0.0	0.0	0.0	0.0	0.88	0.58	
3	0.35	0.0	0.67	0.37	0.0	0.00	0.0	0.0	0.0	0.0	0.96	0.48	
4	0.66	0.0	0.68	0.42	0.0	0.00	0.0	0.0	0.0	0.0	1.04	0.50	

(875, 19)

We use the Pearson correlation to measure the linearity between these features in order to characterize an individual's motion signature. With a total of 19 facial/head features, we compute the Pearson correlation between all 19 of these features, yielding 171 unique pairs of features across all video clips.

The next step involves creating final dataset for our analysis using these correlation values by taking the upper diagonal of the correlation matrix and vectorizing unique values row wise, belonging to each video together, to form a dataframe.

At this stage, the target variable “Real” was added to identify which observation belonged to a fake video and vice versa. The features in this dataset are correlation values obtained from the correlation matrix for each video. Below is the snippet of the final dataset and its dimensions.

	AU01_r	AU02_r	AU04_r	AU05_r	AU06_r	AU07_r	AU09_r	AU10_r	AU12_r	AU14_r	AU15_r	AU17_r	AU20_r	AU25_r	AU26_r	AI
0	0.580848	-0.388830	0.289504	-0.354376	-0.308196	0.142992	-0.390013	-0.283119	-0.384306	0.488931	0.192848	-0.039182	-0.063942	0.067019	-0.1	
1	0.846448	0.160792	0.485862	-0.378533	0.540558	-0.158666	-0.227037	-0.268912	-0.369665	0.263809	0.058232	-0.255885	-0.091503	0.119169	-0.1	
2	0.717958	0.385666	0.427273	0.060000	0.060841	0.148148	0.268729	0.287967	-0.028982	-0.035514	-0.160199	0.000000	0.000000	0.000000	0.000000	-0.1
3	0.767385	0.138281	-0.176998	0.009181	-0.100836	-0.001857	-0.123262	-0.093450	0.425487	-0.097959	-0.206486	0.520681	0.268505	0.000000	0.000000	-0.1
4	0.301615	-0.109730	0.001275	-0.047978	0.161038	0.444690	0.606318	0.222923	0.000000	-0.062383	0.041358	-0.052826	0.000000	0.000000	0.000000	-0.1
5	0.789893	-0.110567	0.335062	-0.482251	-0.399253	-0.140033	-0.398301	-0.490971	-0.265410	0.151105	0.138331	0.082882	-0.302393	-0.002116	-0.1	
6	0.796969	0.079891	0.211918	-0.275550	-0.265850	-0.101404	-0.182434	-0.318659	-0.314976	0.144010	0.377788	0.200915	-0.170320	0.064646	-0.1	
7	0.571886	0.040540	0.212149	-0.745465	-0.616175	-0.193776	-0.703776	-0.081254	0.150844	-0.193589	0.226362	0.373783	0.484608	0.484608	0.000000	-0.1
8	0.623920	-0.16683	0.664321	-0.370781	-0.160165	0.191418	-0.242044	-0.164566	-0.255306	0.066711	0.000000	-0.000000	0.000000	0.000000	0.000000	-0.1
9	0.665425	0.043154	0.070262	-0.585873	-0.425941	-0.170750	-0.478741	-0.452684	-0.423210	0.265037	0.074550	0.099444	-0.357549	-0.094875	-0.1	

## Logistic Regression - Results

Logistic regression approach provides an accuracy of 50% which is equivalent to an equally likely chance of a model detecting a fake video as real and vice versa. Possible justification behind accuracy being ½ can be high dimensionality. In other words, number of observations (N) are less than the number of features /variables (P) which leads to incorporate less samples and encounter an undersampling issue.

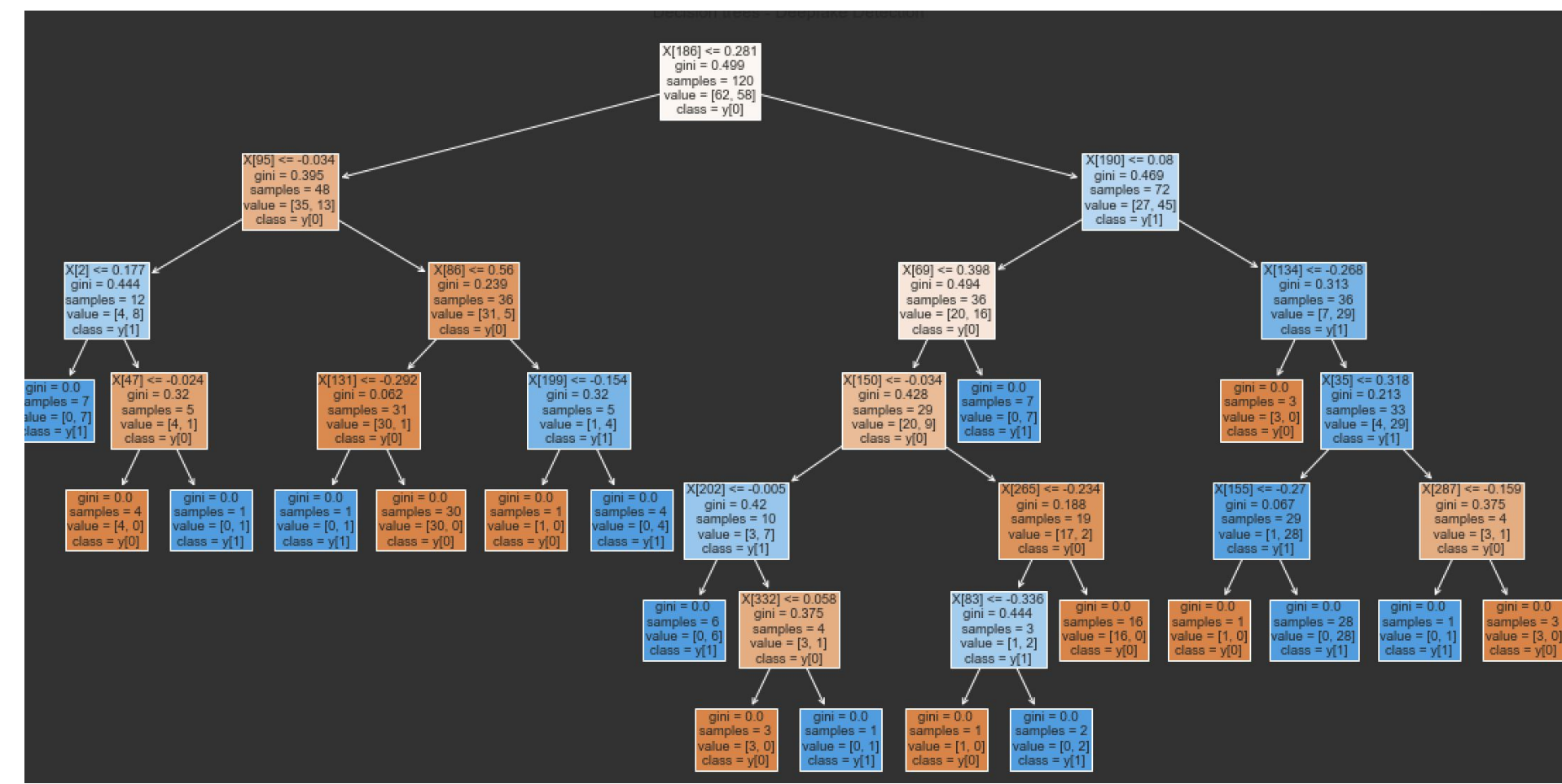
Below is the visual representation of confusion matrix for Logistic Regression.



- Support for this work is greatly acknowledged from American University's School of Arts and Sciences (Data Science Department) and Professor Zois Boukouvalas.
1. Protecting World Leaders Against Deep Fakes (Agarwal, Shruti, Farid, Hany, Gu, Yuming, He, Mingming, Nagano, Koki and Li, Hao), In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops, IEEE, 2019
  2. Nguyen, Thanh & Nguyen, Cuong M. & Nguyen, Tien & Nguyen, Duc & Nahavandi, Saied. (2019). Deep Learning for Deepfakes Creation and Detection: A Survey.
  3. Mal-uses of AI-generated Synthetic Media and Deepfakes: Pragmatic Solutions Discovery Convening, June 2018

## Decision Trees - Results

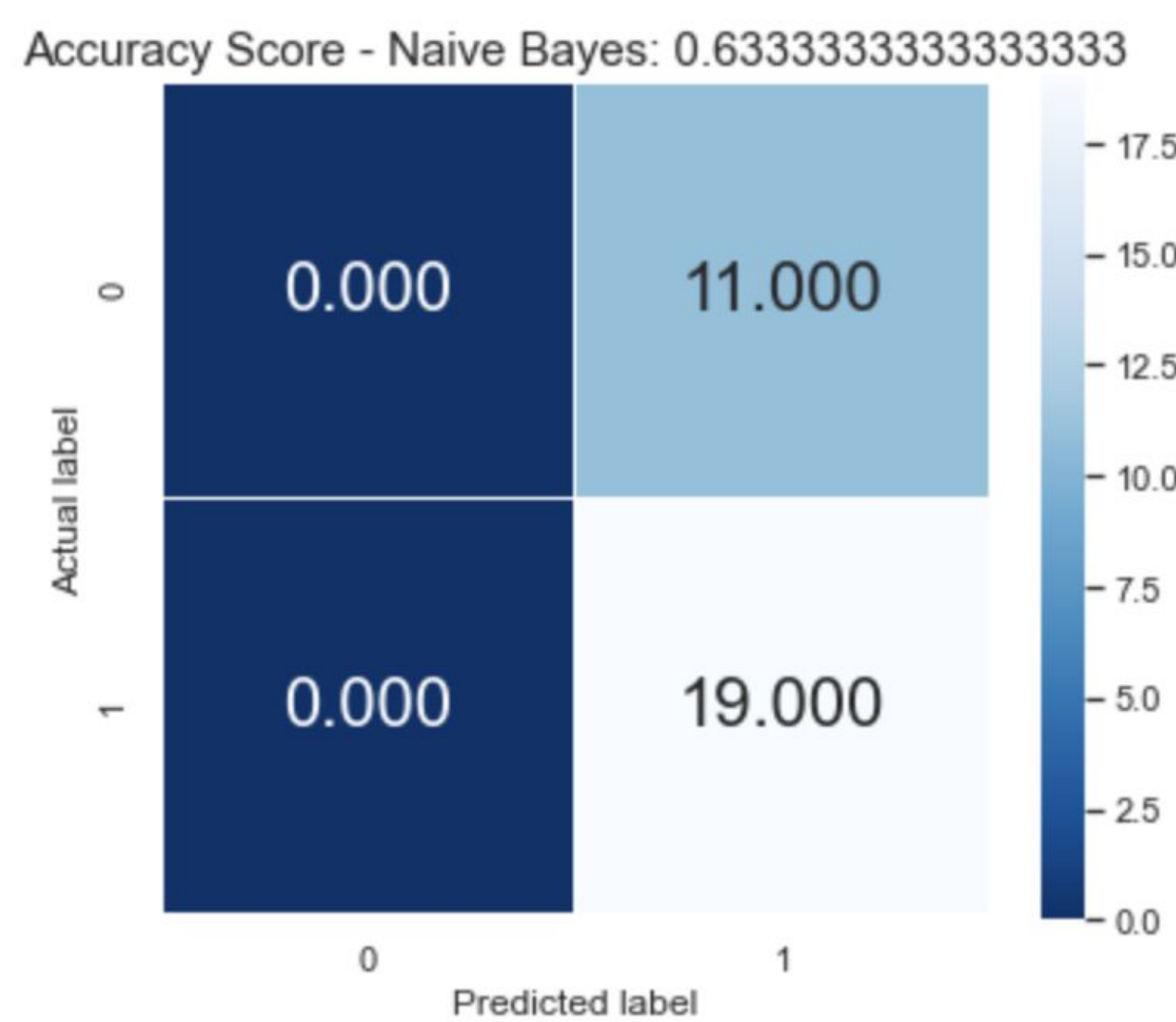
Decision tree algorithm was further utilized to improve accuracy. The percentage increased to 56%. Below is the decision tree created using the gini criterion. Blue box indicates the probability of a video as real and orange box indicates the probability of being fake. Darker shade represents higher probability.



## Naive Bayes - Results

The Naive Bayes classification algorithm provides a higher accuracy than decision trees i.e. it provides 63.3% accuracy. A confusion matrix has been designed in order to visualize the reasoning behind its higher accuracy. It is evident from the matrix that the model is able to perfectly classify the real videos. However, it is totally misclassifying the fake videos as real. Potential issue in this case are in line with issues behind poor classification accuracy provided by logistic regression i.e. undersampling and high dimensionality. Some features with zero vectors should be removed for attaining better accuracy.

Below is the visual representation of confusion matrix for Naive Bayes.



## Conclusion & Future Direction

The Naive Bayes and Decision Trees framework exploits the relationship between independent correlated facial features and our binary target variable in order to detect deep fakes.

Future work will focus on:

- Reducing dimensionality of the data using PCA and Feature Selection Methods
- Using advanced modeling techniques such as Convolutional Neural Network, 2D CNN/ RNN combined and 3D ConvNet