# EDS6344 - Artificial Intelligence for Engineers

## *Hepatitis C Virus (HCV) for Egyptian patients*

## Spring 2024

Professor: Raja Loganantharaj

Group 12

| | |
|---|---|
| **Ashwitha Reddy Nimmala** | **2313782** |
| **Bala Srimani Durga Devi Chikkala** | **2290061** |
| **Azharmadani Syed** | **2316906** |
| **Venkata Sumanth Reddy Vangala** | **2311953** |
| **Amaan Syed** | **2213953** |

# Table of Contents:

# Abstract

This project focuses on developing a classification model for predicting Hepatitis C Virus (HCV) outcomes specifically tailored to Egyptian patients. Through data preprocessing, feature selection, and utilizing various machine learning algorithms, including Logistic Regression, Decision Trees, Random Forest, Gradient Boosting, Support Vector Machines, and Neural Networks, the aim is to accurately predict HCV-related outcomes. Evaluation metrics are used to assess model performance, aiming to improve treatment strategies and patient outcomes in the Egyptian context of HCV infection.

# 1.Introduction

Hepatitis C Virus (HCV) infection poses a significant health challenge in Egypt, with the country experiencing the highest prevalence rate worldwide. In response, a novel Model of Care (MOC) was established in 2006, led by the National Committee for Control of Viral Hepatitis (NCCVH). This review outlines the development and implementation of the Egyptian MOC, which aims to contain the HCV epidemic, provide patient care, and ensure treatment access. The NCCVH's efforts have resulted in over one million patients being evaluated, with more than 850,000 receiving treatment since the program's inception. The establishment of a nationwide network of digitally connected treatment centers has enhanced treatment accessibility. Additionally, localized practice guidelines tailored to local circumstances have been issued and regularly updated. Egypt's successful experience with the MOC serves as a model for other countries and organizations seeking to establish effective programs for HCV care and management.

# 2. DataSet

We collected the dataset from UCI machine learning repository. It consists of 1385 instances and 28 features. The attributes are Age, Gender, BMI(Body Mass Index), Fever, Nausea/Vomiting, Headache, Diarrhea, Fatigue, Bone ache, Jaundice, Epigastria pain, WBC(White Blood Cells), RBC(Red Blood Cells) , HGB(Hemoglobin), Plat(Platelet) , AST1(1 week), ALT1(1 week), ALT4(4 weeks), ALT12(12 weeks), ALT24(24 weeks), ALT36(36 weeks), ALT48(48 weeks), RNA Base, RNA 4, RNA 12, RNA EOT , RNA EF(Elongation Factor), Baseline Histological Grading. The target variable is Baseline Histological Staging.

# 3. Data Preprocessing

## 3.1 Handling Missing values

This is the basic step, where we identify all the missing values in the dataset and fill them with mean, median and mode. In our data set there are no missing values.

Fig 1. Missing Values

## 3.2 Identifying Duplicates

We identified duplicates using.duplicated().sum() function and found that there are no duplicates in the dataset.

## 3.3 Skewness

Skewness is a statistical measure that describes the asymmetry of a probability distribution or a dataset. In our dataset only RNA_12 is right skewed, rest all the other columns are having a normal distribution. There is no need of transformations as well.



Fig 2. Skewness

- Only RNA_12 is right skewed, rest all the other columns are having a normal distribution.
- No need of transformations as well.

### 3.4 Outliers
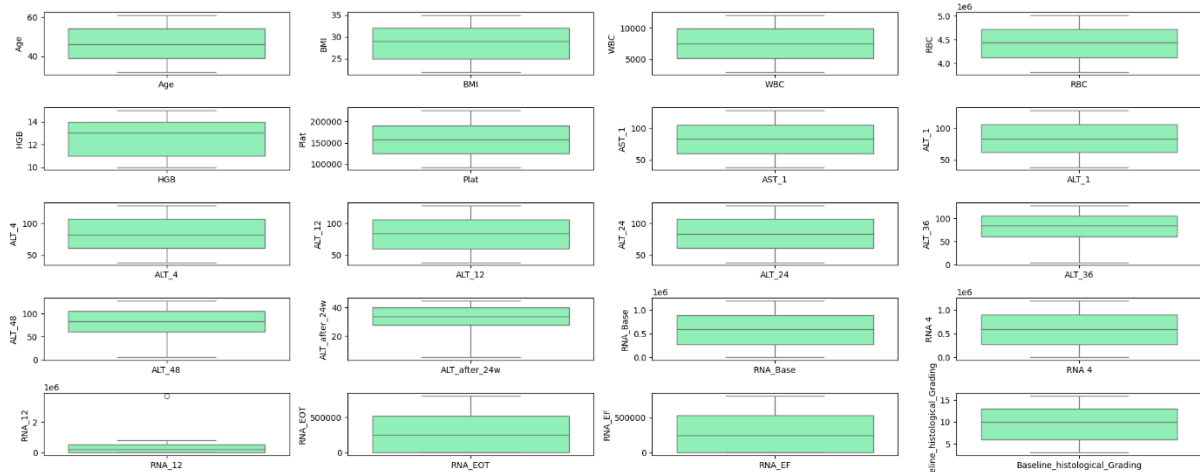
There is only one outlier from the whole dataset.



Fig 3. Outliers

From the above plot we can see that RNA12 has one outlier.

# 4. Exploratory Data Analysis (EDA)

In this section, we visualized the data using various graphs such as heatmaps, bar graphs etc., and analysed the patterns within the data to draw some basic conclusions.



Fig 4. Correlation map.

From the image, it seems that most variables have very low correlation with each other, as indicated by the values close to zero and the light coloration. This lack of strong correlation suggests that there may not be any single pair of variables that strongly predict one another.

There are a few exceptions with moderate correlations. For example, at the end of the columns, you can see a correlation of 0.4 between RNA_EF and Baselinehistological_Grading, suggesting a moderate positive relationship.

## Analysis Of Target Variable

In our dataset the target variable is baseline histological staging. Baseline histological staging typically involves categorizing individuals with HCV infection into different stages based on the severity of liver disease. It consists of different stages like portal fibrosis, few septa, many septa, cirrhosis.



Fig 5. Target Variable distribution.

From the above graph we can see that all the classes are balanced.

## Categorical vs Target analysis

Fig 6. Categorical vs Target variable

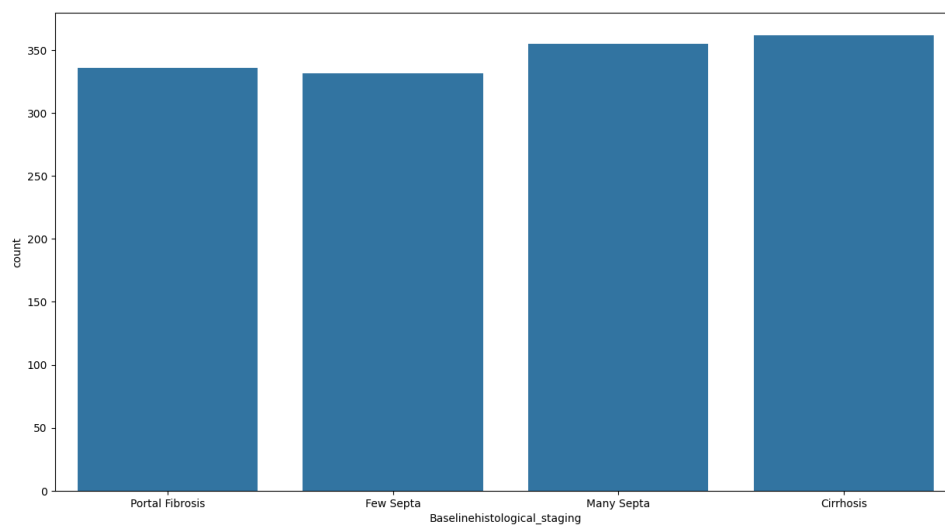This plot is a stacked bar chart that shows the distribution of Baselinehistological Staging of Hepatitis C Virus (HCV) patients based on the presence or absence of epigastric pain.

Each bar represents a category of epigastric pain: "Absent" on the left and "Present" on the right. Within each bar, different colors represent the stages of liver disease:

- Blue represents "Portal Fibrosis"
- Orange represents "Few Septa"
- Green represents "Many Septa"
- Red represents "Cirrhosis"

The height of each color segment within the bars indicates the number of patients in each histological stage with either absent or present epigastric pain.

The chart aims to compare the severity of liver disease in patients with and without epigastric pain. For instance, if the red section (Cirrhosis) is significantly larger in the "Present" category, it could suggest that patients with epigastric pain are more likely to have more advanced liver disease. Conversely, if the blue section (Portal Fibrosis) is higher in the "Absent" category, it might indicate that patients without epigastric pain are more often in the early stages of liver disease.

# Checking histological staging across each gender



Fig 7. Histological Staging.

The chart allows for a comparison between the number of males and females across different stages of liver disease. For instance, it seems that in the "Portal Fibrosis" stage, there are slightly more males than females, while in the "Cirrhosis" stage, the numbers are nearly equal between genders.

This type of visualization helps in understanding gender distribution within each histological stage and can be important for medical research and resource allocation in healthcare settings.

# Symptoms vs Histological stages across Gender

The distribution of histological stages across genders for various symptoms like fever, nausea/vomiting, headache, etc. Parameters like plot order, hue order, height, aspect ratio, and colour palette are customized for each plot to enhance visualization.



Fig 8. Histological Staging.

- On the left side, the chart titled "Gender = Male" shows the count of male patients with and without jaundice across the four stages. Similarly, on the right side, the chart titled "Gender = Female" shows the count of female patients.
- The 'count' on the y-axis likely represents the number of patients, while the 'Jaundice' on the x-axis is divided into 'Present' and 'Absent', indicating whether the patients had jaundice.
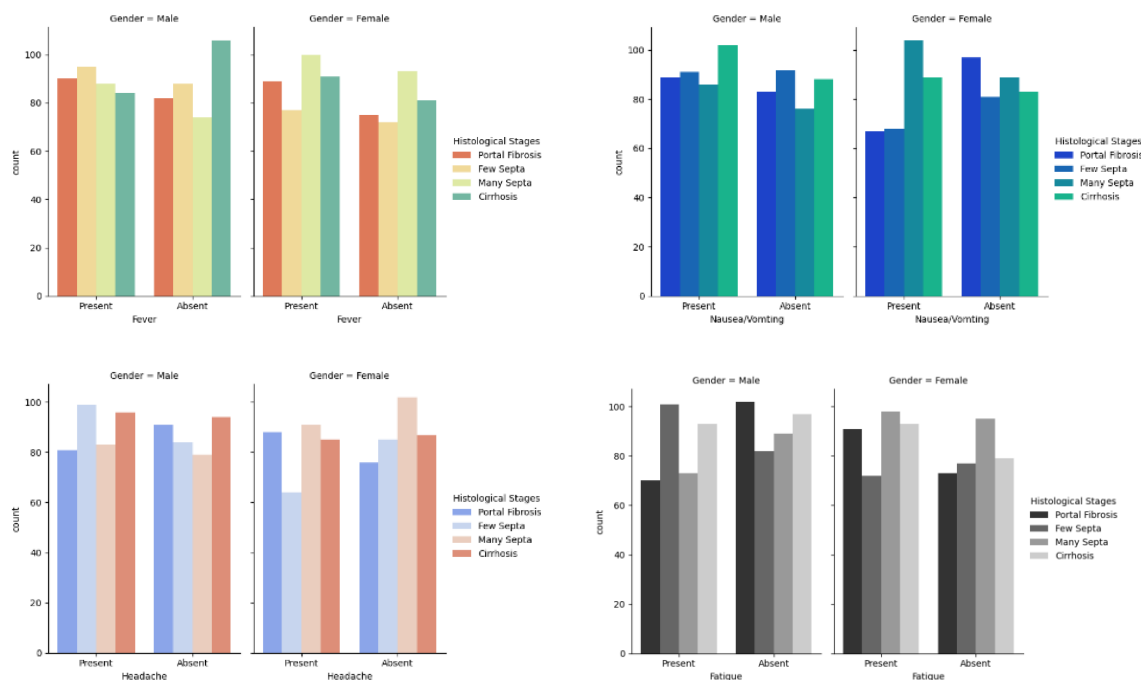- From the charts, we can compare the frequency of jaundice in males and females across different stages of liver disease. It seems that for both genders, jaundice is present across all histological stages, with the counts varying between stages.

## Alanine Transaminase

The DataFrame df_alt contains data related to ALT (Alanine Transaminase) levels at various time points or intervals, including ALT measurements at 1, 4, 12, 24, 36, and 48 weeks, as well as ALT levels after 24 weeks.

| | ALT_1 | ALT_4 | ALT_12 | ALT_24 | ALT_36 | ALT_48 | ALT_after_24w |
|---|---|---|---|---|---|---|---|
| 0 | 84 | 52.0 | 109 | 81 | 5 | 5 | 5 |
| 1 | 123 | 95.0 | 75 | 113 | 57 | 123 | 44 |
| 2 | 49 | 95.0 | 107 | 116 | 5 | 5 | 5 |
| 3 | 64 | 109.0 | 80 | 88 | 48 | 77 | 33 |
| 4 | 104 | 67.0 | 48 | 120 | 94 | 90 | 30 |

Fig 9. Alanine Transaminase.

## Distribution for ALT (per each week) and [stages,grading,HBG] for both gender

The function ALTSwarmPlot creates swarm plots to visualize ALT (Alanine Transaminase) levels across different histological stages and genders. It takes parameters such as the dataset (data), x-axis variable (x), hue variable (hue), figure size (figsize), and color palette (palette). However, there seems to be a discrepancy in the function call regarding the data source and variables used.



Fig 10. ALT Weeks vs Stages.

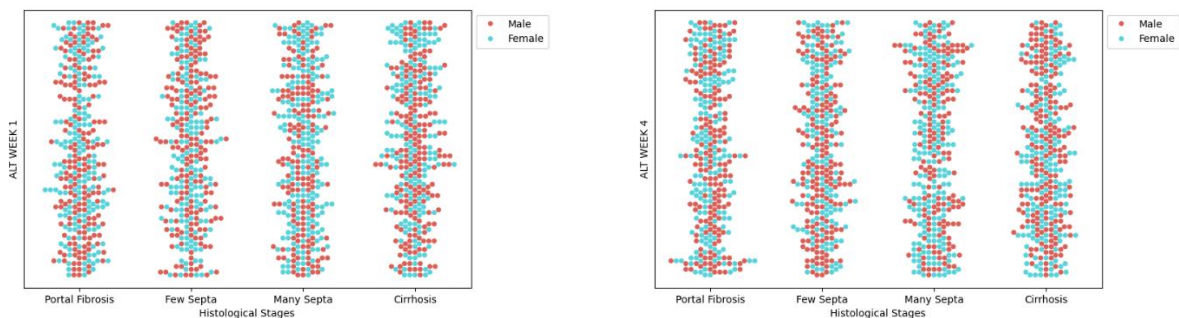Each scatter plot is distributed vertically by ALT levels for a specific week post-treatment, indicating how ALT levels vary over time in relation to the stage of liver disease. The arrangement allows for a comparison of ALT levels at different times during treatment and monitoring, and whether there are differences in these levels between male and female patients within each stage of liver disease.

A few observations that can be inferred from the plots might include:

- ALT levels vary within each stage and over the weeks measured.
- There is a mix of ALT levels between male and female patients at each time point, without any immediately obvious pattern or distinction between genders across the stages.
- By comparing across the different weeks, one could assess whether the ALT levels generally increase, decrease, or remain stable as time progresses, which can be indicative of the progression or improvement of liver condition in response to treatment.

## RNA Stages

The DataFrame RNAdata likely contains data related to RNA measurements. Without seeing the data, I can't provide specifics, but generally, it might include information such as gene expression levels, RNA sequencing data, or other RNA-related metrics. If you provide a glimpse of the first few rows using the head() function, I can give you a more detailed description of its structure and contents.

|   | RNA_Base | RNA 4 | RNA_12 | RNA_EOT | RNA_EF |
|---|---|---|---|---|---|
| 0 | 655330 | 634536 | 288194 | 5 | 5 |
| 1 | 40620 | 538635 | 637056 | 336804 | 31085 |
| 2 | 571148 | 661346 | 5 | 735945 | 558829 |
| 3 | 1041941 | 449939 | 585688 | 744463 | 582301 |
| 4 | 660410 | 738756 | 3731527 | 338946 | 242861 |

Fig 11. RNA Stages.

## RNA VS STAGES

This function call attempts to create a swarm plot to visualize RNA data across different histological stages and genders. It specifies parameters such as the data source (data), the variable for the x-axis (x), the figure size (figsize), the color palette (palette), and the variable for coloring the plot (hue). However, there might be a potential mismatch between the specified data sources (df and df_cat), which could cause issues unless they are appropriately related.

Fig 12. RNA vs Stages.

From the plots, we can see that:

- RNA levels are distributed across all stages of liver disease for both genders.
- At the RNA Base level, RNA 4, and RNA 12, there is a wide distribution of viral load levels among patients in all stages, suggesting variability in how patients present and respond to treatment.
- The RNA EOT and RNA EF scatter plots indicate a decrease in viral load, showing many points clustered at the lower end of the y-axis, which would be expected if patients respond well to treatment.

## Statistical Analysis for Categorical Columns
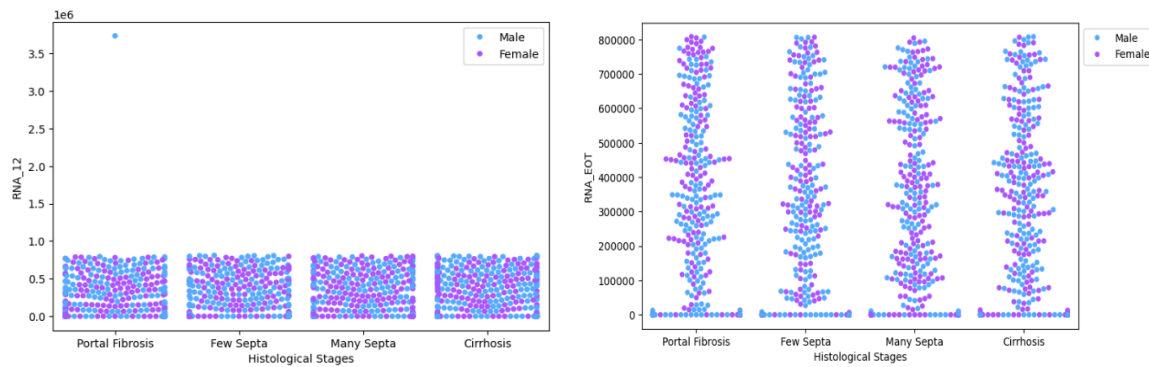
In the Categorical Columns we used a chi-square test for independence for each categorical variable in the DataFrame df_cat, excluding 'Baselinehistological_staging'. It calculates contingency tables for each variable and 'Baselinehistological_staging', then computes chi-square statistics and p-values using chi2_contingency from scipy.stats. The results, indicating the significance of the relationship between each categorical variable and 'Baselinehistological_staging' based on p-values, are printed.

```
Chi-square test for Gender: p-value = 0.08318737457020528
Chi-square test for Fever: p-value = 0.5356332988397037
Chi-square test for Nausea/Vomting: p-value = 0.16217902183179256
Chi-square test for Headache: p-value = 0.9821273061728452
Chi-square test for Fatigue: p-value = 0.589280639706306
Chi-square test for Jaundice: p-value = 0.29896566634958105
Chi-square test for Diarrhea: p-value = 0.7010432291061677
Chi-square test for Epigastric_pain: p-value = 0.08759780229199586
```

Fig 13. Chi-Square Test.

From the chi-square tests, p-values lower than a selected level of significance (e.g., 0.05) usually indicate a significant association between the categorical variable and target variable. In our case no p-values are less than 0.05, which means that non of these factors have significant association with the target variable Baselinehistological_staging at 0.05 level.

However, choosing significance level is arbitrary and could be adjusted according to our specific needs and circumstances of our analysis. Statistically, while significance is critical in evaluating whether or not a predictor is important; there are other things to consider as well.

## Statistical Analysis for Numerical Columns

We applied ANOVA (Analysis of Variance) tests for each numerical variable in a DataFrame (df_num), comparing them against the target variable 'Baselinehistological_staging'. It calculates the F-statistic and associated p-value for each variable, indicating if there are significant differences in means across different categories of the target variable. The results help assess the importance of numerical variables in predicting the target variable.

```
ANOVA for Age: p-value = 0.0
ANOVA for BMI: p-value = 0.0
ANOVA for WBC: p-value = 0.0
ANOVA for RBC: p-value = 0.0
ANOVA for HGB: p-value = 0.0
ANOVA for Plat: p-value = 0.0
ANOVA for AST_1: p-value = 0.0
ANOVA for ALT_1: p-value = 0.0
ANOVA for ALT_4: p-value = 0.0
ANOVA for ALT_12: p-value = 0.0
ANOVA for ALT_24: p-value = 0.0
ANOVA for ALT_36: p-value = 0.0
ANOVA for ALT_48: p-value = 0.0
ANOVA for ALT_after_24w: p-value = 0.0
ANOVA for RNA_Base: p-value = 0.0
ANOVA for RNA 4: p-value = 0.0
ANOVA for RNA_12: p-value = 5.939803316101762e-251
ANOVA for RNA_EOT: p-value = 1.194969937533851e-281
ANOVA for RNA_EF: p-value = 4.098330501007371e-282
ANOVA for Baseline_histological_Grading: p-value = 0.0
```

Fig 14. Anova.

The results of the ANOVA tests indicate that all numerical variables have extremely low p-values (close to zero), suggesting that there are significant differences in means across the categories of the target variable (Baselinehistological_staging).

# 5. Model Building

In this section we will be discussing about the algorithms implemented, hyperparameter tuning techniques used along with the feature selection. After doing EDA and data pre-processing, we will be having the cleaned dataset. Using this clean dataset, we split it into train and test dataset. Here, The train_test_split function from the sklearn.model_selection module is employed for this purpose.

## 5.1 Models

**Random Forest Classifier:** Random Forest Classifier is a machine learning algorithm that combines multiple decision trees to make predictions.

```
Train Score: 1.0
Test Score: 0.7617328519855595
```

Fig 15. RF Accuracy.

To visualizes a decision tree from a Random Forest classifier. It selects one of the trees from the Random Forest, exports it to DOT format, creates a graph from the DOT data using pydotplus, and then displays the decision tree as a PNG image.
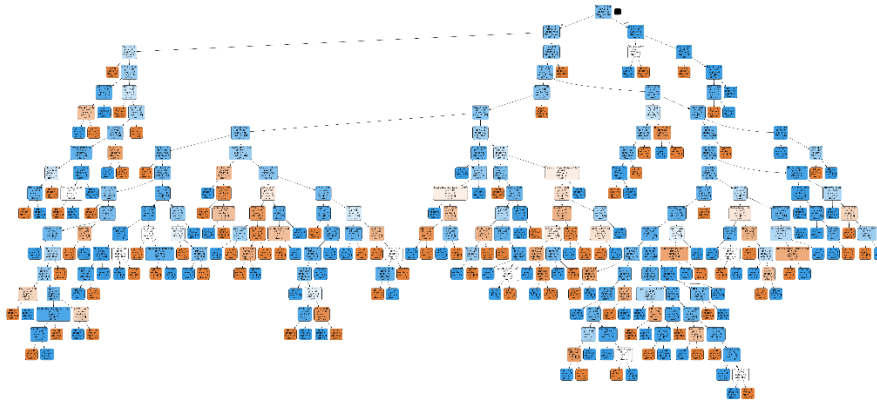
Fig 16. DT from a RF.

**Gradient Boost Classifier:** Gradient Boosting Classifier is an ensemble learning algorithm that iteratively improves predictive accuracy by sequentially adding weak learners to the ensemble, with each one focusing on the errors of its predecessors. It's highly effective for classification tasks, especially with complex datasets, due to its ability to handle non-linear relationships and feature interactions effectively.

```
Train Score: 0.8456678700361011
Test Score: 0.7364620938628159
```

Fig 17. GB Accuracy.

**Decision Tree Classifier:** The Decision Tree Classifier is a machine learning model that makes predictions by recursively splitting the data based on features to create a tree-like structure, resulting in a series of if-else decision rules.

```
Train Score: 1.0
Test Score: 0.5595667870036101
```

Fig 18. DT Accuracy.

**AdaBoost Classifier:** AdaBoost is a machine learning algorithm that combines weak learners sequentially, focusing more on misclassified instances in each iteration to create a strong classifier.

```
Train Score: 0.7680505415162455
Test Score: 0.740072202166065
```

Fig 19. AdaBoost Accuracy.

**Bagging Classifier:** Bagging is a machine learning technique that combines multiple models trained on random subsets of the training data to improve predictive performance. It reduces variance and helps prevent overfitting by averaging the predictions of the individual models, resulting in a more robust and accurate classifier.

```
Train Score: 0.9918772563176895
Test Score: 0.7184115523465704
```

Fig 20. Bagging Accuracy.

**XGB Classifier**: XGBoost classifier is like a team captain making decisions. It trains a series of decision-makers (trees) one after the other. Each new decision-maker learns from the mistakes of the previous ones, gradually improving the team's overall performance. It's efficient, scales well, and is really good at figuring out patterns in data to make accurate predictions.

11

```
Train Score: 1.0
Test Score: 0.7617328519855595
```

Fig 21. XGBoost Accuracy.

**Voting Classifier:** The Voting Classifier combines predictions from multiple individual classifiers to make a final decision. It's a simple yet effective ensemble method that often improves prediction accuracy by considering diverse perspectives from different classifiers.

```
Train Score: 1.0
Test Score: 0.7509025270758123
```

Fig 22. Voting Accuracy.

Here is the scorecard for the above models:

| | Model | Precision | Recall | F1_score | Train_Accuracy | Test_Accuracy |
|---|---|---|---|---|---|---|
| 0 | Random Forest | 0.669453 | 0.722022 | 0.685772 | 1.000000 | 0.722022 |
| 1 | Decision tree | 0.647106 | 0.570397 | 0.598189 | 1.000000 | 0.570397 |
| 2 | Gradient Boost | 0.634849 | 0.678700 | 0.653105 | 0.895585 | 0.678700 |
| 3 | XGB Classifier | 0.633630 | 0.693141 | 0.656548 | 1.000000 | 0.693141 |
| 4 | Ada Boost Forest | 0.634616 | 0.613718 | 0.623406 | 0.757757 | 0.613718 |
| 5 | Voting Classifier | 0.643626 | 0.707581 | 0.666023 | 1.000000 | 0.707581 |
| 6 | Bagging Classifier | 0.631788 | 0.559567 | 0.586999 | 0.996420 | 0.559567 |

Fig 23. Scorecard.

From the scorecard, we can observe the following:

- The **Random Forest** model has a moderate recall but the lowest precision and a considerable gap between training and test accuracy, indicating potential overfitting.

- The **Decision Tree** model has the lowest recall and test accuracy, suggesting that it might be too simple to capture the complexity of the data.

- The **Gradient Boost** model shows better balance between precision and recall, with a higher F1 score than Random Forest and Decision Tree, and a good balance between train and test accuracy.

- The **XGB Classifier** has the highest F1 score, indicating a good balance between precision and recall. However, it also displays a perfect training accuracy, which could be a sign of overfitting, as reflected in the test accuracy.

- **Ada Boost Forest** shows good recall and a decent F1 score, with more balanced accuracy scores between training and testing, indicating a better generalization than some of the other models.

- The **Voting Classifier** appears to have a good balance of precision, recall, and F1 score, with a high training and test accuracy, which might make it the best generalized model among those listed.

- Lastly, the **Bagging Classifier** has the highest precision and a strong F1 score, with very high training accuracy and reasonably high test accuracy, which could also indicate good generalization but with a slight inclination towards overfitting.

The decision on which model to use will depend on the specific use case and whether precision or recall is more important. For instance, in medical diagnostics, high recall might be preferred to ensure as few false negatives as possible, even at the expense of precision.

# ROC_AUC Curves

This function creates subplots to display ROC curves for multiple models. It calculates the true positive rate (sensitivity) against the false positive rate (1-specificity) for each model's predictions. It then plots the ROC curves with corresponding AUC (Area Under the Curve) values, comparing the performance of different models in distinguishing between classes.
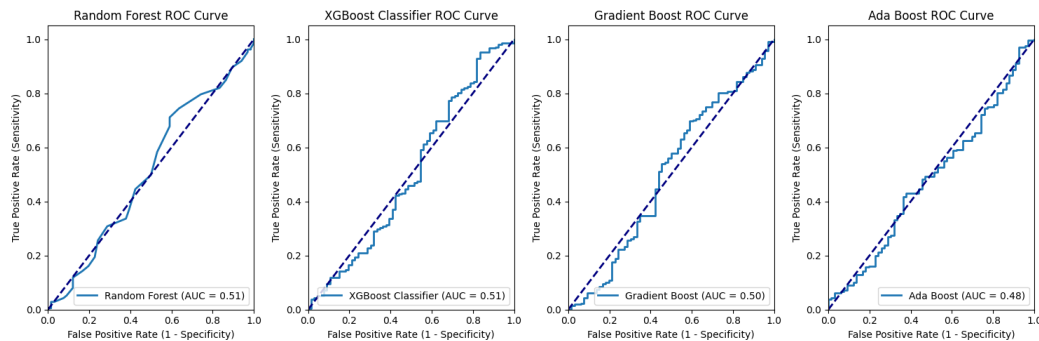


Fig 24. Roc_Auc Curves.

From the curves we can indicate that:

- **Random Forest ROC Curve:** Shows the performance of the Random Forest classifier with an Area Under Curve (AUC) of 0.51. An AUC of 0.5 indicates that the model has no discrimination capability to distinguish between positive and negative classes.
- **XGBoost Classifier ROC Curve:** Shows the performance of the XGBoost classifier with an AUC of 0.51, which is similar to flipping a coin.
- **Gradient Boost ROC Curve:** Shows the performance of the Gradient Boost classifier with an AUC of 0.49. This is slightly less than 0.5, which suggests that the model is performing worse than random guessing.
- **Ada Boost ROC Curve:** Shows the performance of the Ada Boost classifier with an AUC of 0.48, also indicating a performance worse than random guessing.

In general, an AUC closer to 1 indicates a better-performing model, while an AUC closer to 0.5 suggests no predictive power. These particular models are not performing well on the dataset since their AUC scores are near or below 0.5, suggesting that they have no better accuracy than random chance.

# Confusion Matrix

This function generates and plots confusion matrices, along with classification reports, for multiple models. It creates subplots to display each model's results, including predicted versus actual labels. The function allows for easy comparison of model performance in terms of classification accuracy, precision, recall, and F1-score.



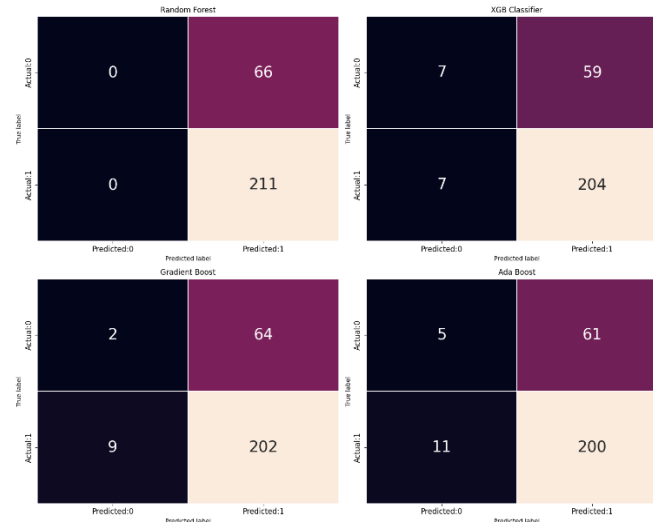Fig 25. Baselinehistological_staging count.

Fig 26. Confusion Matrix.

- Very interesting point to be noticed from this confusion matrix.
- We have very low scores for class '0' label. Why?
- The reason is there is huge imbalance in the data. There are only 336 '0' classes and 1049 '1' classes. So, the model is predicting everything as class '1'.
- We can apply SMOTE to handle this imbalance in the data.

## SMOTE

SMOTE stands for Synthetic Minority Over-sampling Technique. It's a technique used in machine learning to address class imbalance problems, particularly in classification tasks where one class (the minority class) is significantly underrepresented compared to the other class (the majority class).

From the above confusion matrix, we can say that there is huge imbalance in the data. There are only 336 '0' classes and 1049 '1' classes. So, the model is predicting everything as class '1'. We can apply SMOTE to handle this imbalance in the data.

Here is the score card after applying smote.

| | Model | Precision | Recall | F1_score | Train_Accuracy | Test_Accuracy |
|---|---|---|---|---|---|---|
| 0 | Random Forest | 0.669453 | 0.722022 | 0.685772 | 1.000000 | 0.722022 |
| 1 | Decision tree | 0.647106 | 0.570397 | 0.598189 | 1.000000 | 0.570397 |
| 2 | Gradient Boost | 0.634849 | 0.678700 | 0.653105 | 0.895585 | 0.678700 |
| 3 | XGB Classifier | 0.633630 | 0.693141 | 0.656548 | 1.000000 | 0.693141 |
| 4 | Ada Boost Forest | 0.634616 | 0.613718 | 0.623406 | 0.757757 | 0.613718 |
| 5 | Voting Classifier | 0.632843 | 0.700361 | 0.657777 | 1.000000 | 0.700361 |
| 6 | Bagging Classifier | 0.631788 | 0.559567 | 0.586999 | 0.996420 | 0.559567 |

Fig 27. Smote Scorecard.

From the scorecard, we can infer the following:
- Models like Random Forest, XGB Classifier, and Voting Classifier have perfect training accuracy, indicating they may be overfitting the training data, as reflected by lower test accuracy.

14

- The Gradient Boost model has the highest test accuracy and a relatively high F1 score, suggesting it might be the best at generalizing to new data among the listed models.
- Precision and recall are more balanced in models such as Random Forest and Gradient Boost, which could be important if the cost of false positives and false negatives is similar.
- The Decision Tree model has lower test accuracy compared to others, indicating that it may not perform as well on unseen data.
- The F1 scores are moderate, which indicates that there is room for improvement in the models' ability to balance precision and recall.

## Smote ROC_AUC Curve

The function creates ROC curves for multiple machine learning models, such as Random Forest, XGBoost Classifier, Gradient Boost, and Ada Boost. It uses the plot_roc_curves_subplots function to display the ROC curves and their corresponding AUC values in subplots for each model. The models and test data are provided as inputs to evaluate the models' performance.



Fig 28. Smote Roc Curve.

In this case, all models have an AUC of less than 0.5, which is unusual and suggests that the models may be performing worse than random guessing. This can happen when the models are predicting inversely - that is, predicting a negative when it's actually positive, and vice versa. It may also be an indication that the models are not well-configured for the dataset or that the data itself may not contain clear signals to differentiate between the classes effectively.

## Feature Selection

Feature selection is essential step for building accurate and effective models. It helps in avoiding the overfitting, multicollinearity and enhances the capability of the models.



Fig 29. Feature Selection.

After using feature selection, there is no improvement in the scores

| | Model | Precision | Recall | F1_score | Train_Accuracy | Test_Accuracy |
|---|---|---|---|---|---|---|
| 0 | Random Forest | 0.611979 | 0.664260 | 0.634427 | 1.000000 | 0.664260 |
| 1 | Decision tree | 0.621306 | 0.592058 | 0.605354 | 1.000000 | 0.592058 |
| 2 | Gradient Boost | 0.635344 | 0.671480 | 0.650938 | 0.902745 | 0.671480 |
| 3 | XGB Classifier | 0.673451 | 0.722022 | 0.688873 | 1.000000 | 0.722022 |
| 4 | Ada Boost Forest | 0.592595 | 0.534296 | 0.559056 | 0.754177 | 0.534296 |
| 5 | Bagging Classifier | 0.610275 | 0.566787 | 0.585817 | 0.997017 | 0.566787 |

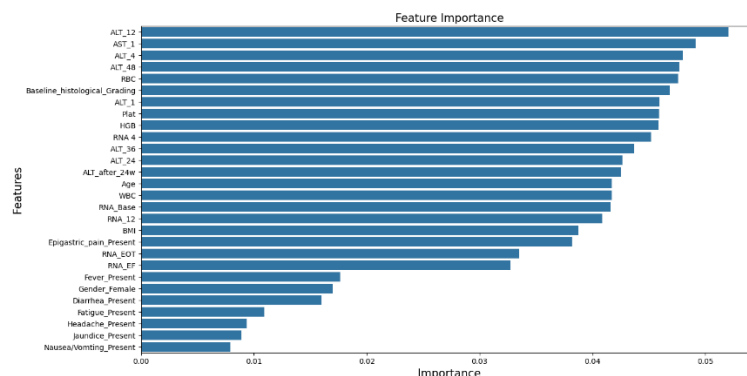Fig 30. Feature Selection Scorecard.

After applying feature selection XGB Classifier excels with improved precision, recall, and F1 score, maintaining consistent test accuracy. Gradient Boost and Random Forest show steady performance, while Ada Boost Forest declines. Decision Tree and Bagging Classifier perform moderately, with slight improvements. Overall, XGB Classifier remains the top performer.

- Even after selecting the best parameters and building the model, we are unable to increase the scores for class label 0.
- Now what we can do is applying PCA on the data.
- After that lets try to tune the model using GridSearchCv.

## 5.2 Principal Component Analysis:

Principal Component Analysis, is a statistical technique used to reduce the dimensionality of high-dimensional data while retaining most of its variability.

Here are the results after applying PCA to the data. We can observe that there is an increase in the score.

| | Model | Precision | Recall | F1_score | Train_Accuracy | Test_Accuracy |
|---|---|---|---|---|---|---|
| 0 | Random Forest | 0.822109 | 0.821429 | 0.821302 | 1.000000 | 0.821429 |
| 1 | Decision tree | 0.663892 | 0.663690 | 0.663503 | 1.000000 | 0.663690 |
| 2 | Gradient Boost | 0.707871 | 0.702381 | 0.700149 | 0.911194 | 0.702381 |
| 3 | XGB Classifier | 0.805598 | 0.803571 | 0.803181 | 1.000000 | 0.803571 |
| 4 | Ada Boost Forest | 0.648844 | 0.648810 | 0.648735 | 0.741791 | 0.648810 |
| 5 | Bagging Classifier | 0.765324 | 0.755952 | 0.753576 | 0.991045 | 0.755952 |

Fig 31. PCA Scorecard.

- The Random Forest and XGB Classifier show a perfect training accuracy of 1.000000, indicating they may have overfitted the training data. Their test accuracies are good but a perfect score on training data warrants a closer look to confirm the models' ability to generalize.
- The Decision Tree and Gradient Boost models show reasonable performance, though not as high as Random Forest or XGB Classifier. Ada Boost Forest has the lowest precision, recall, F1 score, and test accuracy, indicating it may be the least effective model after applying PCA in this scenario.

- Bagging Classifier presents a good balance between precision and recall and has a high training accuracy but without the perfect fit seen in Random Forest and XGB Classifier, suggesting better generalization.

Overall, applying PCA seems to have helped some models more than others. The Random Forest and XGB Classifier might be considered the best performing models based on these results, but the potential overfitting indicated by their perfect training accuracy could be a point of concern.

## After applying PCA:



Fig 32. PCA Confusion Matrix.

- After applying PCA, the '0' class labels are predicted correctly.
- There are few misclassifications in the dataset but still the model is performing good.
- We have achieved an accuracy of 84% for Random Forest which is the highest.
- Models like Random Forest and XGBoost Classifier show relatively high true positive and true negative rates, suggesting good classification performance.
- The Gradient Boost and Ada Boost models have lower true positive rates compared to the Random Forest and XGBoost Classifier.
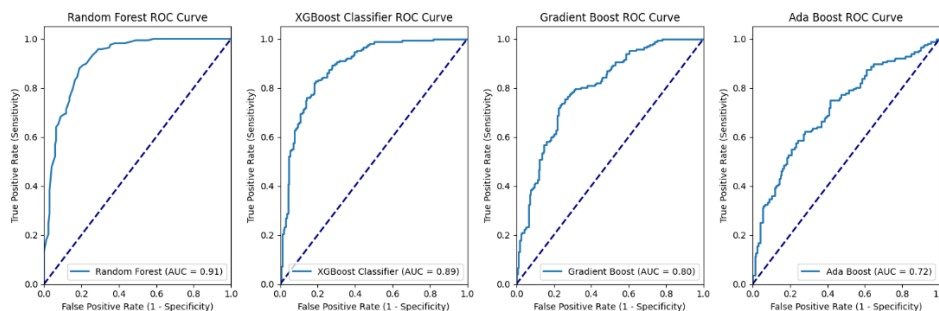


Fig 33. PCA Roc Curve.

- The Random Forest and XGBoost Classifier models have high AUC scores (0.91 and 0.90, respectively), indicating excellent predictive performance.

- The Gradient Boost's AUC is 0.76, which is decent, showing good classification capability.
- The Ada Boost model has the lowest AUC score (0.66), which is still fair but suggests the model's predictive performance is not as strong as the others.

## 5.3 HYPER PARAMETER TUNING:

Hyperparameter tuning is the optimization of a model's configuration settings, called hyperparameters, to enhance its performance. It involves finding the best combination of hyperparameters to maximize the model's accuracy or other performance metrics on a validation dataset. Techniques like grid search or random search are commonly used to search through the hyperparameter space and improve the model's generalization and predictive accuracy.

**Tuning on all the features:**

| | Model | Precision | Recall | F1_score | Train_Accuracy | Test_Accuracy |
|---|---|---|---|---|---|---|
| 0 | Tuned Random Forest | 0.656993 | 0.646209 | 0.651349 | 0.699160 | 0.651349 |
| 1 | Tuned Decision Tree | 0.671536 | 0.559567 | 0.591884 | 0.604577 | 0.591884 |
| 2 | Tuned Bagging Classifier | 0.618112 | 0.469314 | 0.507547 | 0.572746 | 0.507547 |
| 3 | Tuned AdaBoost Classifier | 0.646748 | 0.624549 | 0.634698 | 0.797732 | 0.634698 |
| 4 | Tuned GradientBoost Classifier | 0.617967 | 0.646209 | 0.630902 | 1.000000 | 0.630902 |

Fig 34. HPT Scorecard on all features.

- The Tuned Decision Tree model has the highest test accuracy and also shows a good balance between precision and recall.
- The Tuned Random Forest and Tuned AdaBoost Classifier show relatively high F1 scores, indicating a good balance between precision and recall. However, their test accuracy is slightly lower than the Decision Tree model.
- The Tuned Bagging Classifier has the lowest test accuracy and F1 score, suggesting it might be the least effective of the tuned models.
- The Tuned GradientBoost Classifier shows a very high training accuracy but a much lower test accuracy, indicating potential overfitting.
- The test accuracies are fairly consistent, but not as high as one might expect from tuned models, indicating there may still be room for improvement, either in further tuning or in feature engineering.

**Tuning only on the Important Features:**

| | Model | Precision | Recall | F1_score | Train_Accuracy | Test_Accuracy |
|---|---|---|---|---|---|---|
| 0 | Tuned Random Forest | 0.631319 | 0.527076 | 0.561186 | 0.679741 | 0.561186 |
| 1 | Tuned Decision Tree | 0.642098 | 0.494585 | 0.531148 | 0.597037 | 0.531148 |
| 2 | Tuned Bagging Classifier | 0.631594 | 0.404332 | 0.427047 | 0.513989 | 0.427047 |
| 3 | Tuned AdaBoost Classifier | 0.648105 | 0.642599 | 0.645291 | 0.805737 | 0.645291 |
| 4 | Tuned GradientBoost Classifier | 0.635789 | 0.664260 | 0.648500 | 0.995823 | 0.648500 |

Fig 35. HPT Scorecard on Important features.

- The Tuned AdaBoost Classifier and Tuned GradientBoost Classifier seem to be performing the best in terms of balance across precision, recall, and F1_score, which is an indicator of model performance that takes both precision and recall into account.
- The Tuned Random Forest and Tuned Decision Tree have moderate precision but lower recall, which suggests they are more conservative with positive predictions but may miss a significant number of positive cases.
- The Tuned Bagging Classifier has the lowest F1 score and test accuracy, indicating that it might be the least effective of the tuned models on the important features.
- The Tuned GradientBoost Classifier shows a high level of overfitting with nearly perfect training accuracy but lower test accuracy.

The model with the highest test accuracy after tuning is the GradientBoost Classifier, followed closely by the AdaBoost Classifier. These would likely be the best candidates for further optimization and potentially for use in making predictions, depending on the cost of false positives and false negatives in the specific application context.

## 5.4 Deep Learning

## MLP Classifier on complete data

An MLP Classifier, a type of neural network model, on a given dataset. It evaluates the model's performance by calculating the accuracy on both the training and test datasets. Finally, it prints the train and test accuracies to assess the model's performance in classification tasks.

```
Train Accuracy: 100.00%
Test Accuracy: 68.59%
```

Fig 36. MLP Accuracy.

## MLP Classifier on PCA data
An MLPClassifier neural network model on data that has undergone Principal Component Analysis (PCA) transformation. It evaluates the model's performance on both the training and test datasets and prints the corresponding accuracies. This allows for an assessment of how well the model generalizes to unseen data after dimensionality reduction with PCA.

```
Train Accuracy: 100.00%
Test Accuracy: 80.95%
```
Fig 37. MLP PCA Accuracy.

The classification report suggests that the MLP classifier is performing well on both classes after PCA with good precision, recall, and F1-scores. It also indicates a balanced dataset with almost equal support for both classes. The model shows good generalization ability, as indicated by the high accuracy and balanced metrics across both classes.

## Fully Connected Layer on complete data

The fully connected segment builds, trains, and evaluates a neural network model using TensorFlow's Keras API. It defines a sequential model with dense layers, compiles it with an optimizer and loss function, trains it on training data with specified epochs and batch size, and evaluates its performance on test data, providing the loss value and accuracy.

```
Epoch 97/100
28/28 [==============================] – 0s 4ms/step – loss: 0.0000e+00 – accuracy: 0.7494 – val_loss: 0.0000e+00 –
val_accuracy: 0.7838
Epoch 98/100
28/28 [==============================] – 0s 748us/step – loss: 0.0000e+00 – accuracy: 0.7494 – val_loss: 0.0000e+00
– val_accuracy: 0.7838
Epoch 99/100
28/28 [==============================] – 0s 750us/step – loss: 0.0000e+00 – accuracy: 0.7494 – val_loss: 0.0000e+00
– val_accuracy: 0.7838
Epoch 100/100
28/28 [==============================] – 0s 705us/step – loss: 0.0000e+00 – accuracy: 0.7494 – val_loss: 0.0000e+00
– val_accuracy: 0.7838
9/9 [==============================] – 0s 397us/step – loss: 0.0000e+00 – accuracy: 0.7617

[0.0, 0.7617328763008118]
```

Fig 38. Fully Connected Layer on Complete data Accuracy.

## Fully Connected layer on PCA data

The fully connected segment builds, trains, and evaluates a neural network model using TensorFlow's Keras API. It defines a sequential model with dense layers, compiles it with an optimizer and loss function, trains it on PCA-transformed training data, and evaluates its performance on PCA-transformed test data, providing the loss value and accuracy.

```
Epoch 98/100
34/34 [==============================] – 0s 628us/step – loss: 0.0000e+00 – accuracy: 0.5149 – val_loss: 0.0000e+00
– val_accuracy: 0.4440
Epoch 99/100
click to unscroll output; double click to hide ============] – 0s 634us/step – loss: 0.0000e+00 – accuracy: 0.5149 – val_loss: 0.0000e+00
– val_accuracy: 0.4440
Epoch 100/100
34/34 [==============================] – 0s 642us/step – loss: 0.0000e+00 – accuracy: 0.5149 – val_loss: 0.0000e+00
– val_accuracy: 0.4440
11/11 [==============================] – 0s 348us/step – loss: 0.0000e+00 – accuracy: 0.4970

[0.0, 0.4970238208770752]
```

Fig 39. Fully Connected Layer on PCA data Accuracy.

## CNN layer on complete data

A convolutional neural network (CNN) model using TensorFlow's Keras API. The model comprises 1D convolutional layers followed by max-pooling layers, followed by fully connected layers. It is trained and evaluated on input data (X_train, y_train) and (X_test, y_test) for a specified number of epochs and batch size, providing loss and accuracy metrics.

```
val_accuracy: 0.7838
Epoch 99/100
28/28 [==============================] – 0s 2ms/step – loss: 0.0000e+00 – accuracy: 0.7494 – val_loss: 0.0000e+00 –
val_accuracy: 0.7838
Epoch 100/100
28/28 [==============================] – 0s 2ms/step – loss: 0.0000e+00 – accuracy: 0.7494 – val_loss: 0.0000e+00 –
val_accuracy: 0.7838
9/9 [==============================] – 0s 622us/step – loss: 0.0000e+00 – accuracy: 0.7617

[0.0, 0.7617328763008118]
```

Fig 40. CNN Layer on Complete data Accuracy.

## CNN Model on PCA data

A convolutional neural network (CNN) model using TensorFlow's Keras API on PCA-transformed data. The model consists of convolutional and pooling layers followed by fully connected layers. It's compiled with the Adam optimizer and categorical cross-entropy loss function, trained on the PCA-transformed training data, and evaluated on the PCA-transformed test data, providing loss and accuracy metrics.

```
Epoch 98/100
34/34 [==============================] – 0s 1ms/step – loss: 0.0000e+00 – accuracy: 0.5149 – val_loss: 0.0000e+00 –
val_accuracy: 0.4440
Epoch 99/100
34/34 [==============================] – 0s 1ms/step – loss: 0.0000e+00 – accuracy: 0.5149 – val_loss: 0.0000e+00 –
val_accuracy: 0.4440
Epoch 100/100
34/34 [==============================] – 0s 1ms/step – loss: 0.0000e+00 – accuracy: 0.5149 – val_loss: 0.0000e+00 –
val_accuracy: 0.4440
11/11 [==============================] – 0s 475us/step – loss: 0.0000e+00 – accuracy: 0.4970

[0.0, 0.4970238208770752]
```

Fig 41. CNN Layer on PCA data Accuracy.

- Fully connected layer and CNN layer are not working well on this dataset.
- The accuracy is very less (around 76%).
- Even on the PCA data, Neural Networks is not working well.
- MLP Classifier is working well on the PCA data.

# 6. Conclusion

Based on the work done including various scorecards, confusion matrices, and classification reports, here is a summary and conclusion to help decide the best model:

1. **Random Forest on PCA Data:**
- High precision and recall.
- The ROC curve indicated a very good AUC, and the confusion matrix showed a balanced classification ability for both classes.
- The classification report confirms its good performance, with high F1-scores for both classes.
2. **MLP Classifier on PCA Data:**

- The classification report showed high metrics across the board, with slightly higher precision for class 1 and better recall for class 0 compared to Random Forest.
- However, given the nature of neural networks, they might require more data and careful tuning to generalize well in comparison to ensemble methods like Random Forest.

Considering all metrics and the nature of the data after PCA, We would lean towards the Random Forest model as the best choice for a few reasons:
- It exhibited a strong balance across precision, recall, and F1-score, which is indicative of its robustness.
- They are easier to interpret and diagnose, thanks to the straightforward nature of decision trees that make up the forest, which can be a significant advantage depending on the application area.
- The Random Forest model's performance was consistent across multiple evaluations.

**Conclusion:** In the process of model selection, it's crucial to consider not only the raw performance metrics but also the characteristics of the dataset, the complexity of the models, and the specific context of the problem. For instance, if interpretability is a key concern, ensemble methods like Random Forest would be preferable. If the dataset was larger and more complex, and model interpretability was less of a concern, a neural network approach might be more suitable.

The PCA-transformed dataset appears to be a good fit for the Random Forest model, leading to high-performance metrics while likely offering a good trade-off between accuracy and model complexity. However, if computational resources and data availability allow for it, further exploration with neural networks like the MLP could be beneficial, especially if the dataset's complexity grows or if the performance metrics suggest a non-linear decision boundary that an MLP could navigate better.