

# DATA ANALYSIS USING PYTHON



A Technical project Report  
in partial fulfilment of the degree

**BACHELOR OF TECHNOLOGY**

in

**COMPUTER SCIENCE & ARTIFICIAL INTELLIGENCE**

By

**SYED. UZMA YASMIN**

**2203A52055**

Under the guidance of

**Dr. Ramesh Dadi**

**Assistant Professor, School of CS&AI**

Submitted to



**COMPUTER SCIENCE & ARTIFICIAL INTELLIGENCE**

**SR UNIVERSITY, ANANTHASAGAR, WARANGAL**

**APRIL 2025.**

# 1. HEART FAILURE PREDICTION

## 1. Introduction

Heart failure is a critical global health issue affecting millions of individuals each year. It is a condition where the heart cannot pump blood effectively, leading to inadequate blood flow to meet the body's needs. Early prediction and diagnosis of heart failure can significantly reduce mortality and enhance patient outcomes through timely intervention.

With advancements in data science, machine learning has emerged as a powerful tool for the prediction of diseases, including heart conditions. This report outlines a detailed machine learning pipeline using Python to predict heart failure, leveraging the Kaggle dataset titled "Heart Failure Prediction" by Federico Soriano.

## 2. Dataset Overview

The dataset used in this project is publicly available on Kaggle and contains 918 samples with 12 features each. These features include various medical parameters that are known risk indicators for heart failure. The target variable is heart disease, a binary indicator denoting the presence (1) or absence (0) of heart disease.

### 2.1 Features Description

- **Age:** Age of the person
- **Sex:** Gender (1 = Male, 0 = Female)
- **Chest Pain Type:** Type of chest pain experienced
- **Resting BP:** Resting blood pressure
- **Cholesterol:** Serum cholesterol (mg/dl)
- **Fasting BS:** Fasting blood sugar > 120 mg/dl (1 = True, 0 = False)
- **Resting ECG:** Resting electrocardiogram results
- **Max HR:** Maximum heart rate achieved
- **Exercise Angina:** Exercise-induced angina (1 = Yes, 0 = No)
- **Old peak:** ST depression induced by exercise relative to rest
- **ST\_Slope:** The slope of the peak exercise ST segment
- **Heart Disease:** Target variable (1 = HeartDisease, 0 = No Disease)

### 3. Data Preprocessing and Data Cleaning

Checked for missing values and confirmed the dataset is complete.

- Verified data types for each column.
- **Handling Missing Values:** No missing values were present.
- **Categorical Encoding:** One-hot encoding was applied to categorical features like ChestPainType, RestingECG, ST\_Slope.
- **Normalization:** Z-score normalization was applied to numerical features.
- **Outlier Detection:** Box plots and kurtosis values helped identify outliers.

### 4. Exploratory Data Analysis (EDA)

EDA was performed to understand the distributions and relationships within the data.

- **Age Distribution:** Most patients were between 40–60 years.
- **Heart Disease vs. Sex:** Males were more frequently diagnosed with heart disease.
- **Chest Pain Type:** Atypical angina was a common symptom among those with heart disease.
- **Correlation Matrix:** Old peak, MaxHR, and Exercise Angina showed strong correlations with HeartDisease.

### 5. Model Building and Models Used

- PCA
- Histogram
- Box Plot and Scatter Plot
- Cleaning and Processing
- Logistic Regression
- Evaluation using Confusion Matrix
- K-Nearest Neighbors (KNN)
- Support Vector Machine (SVM)
- Decision Tree
- Naïve Bayes
- Confusion Matrix

- ROC Curve
- Precision-Recall Curve
- Kurtosis

Each model was trained and evaluated using accuracy.

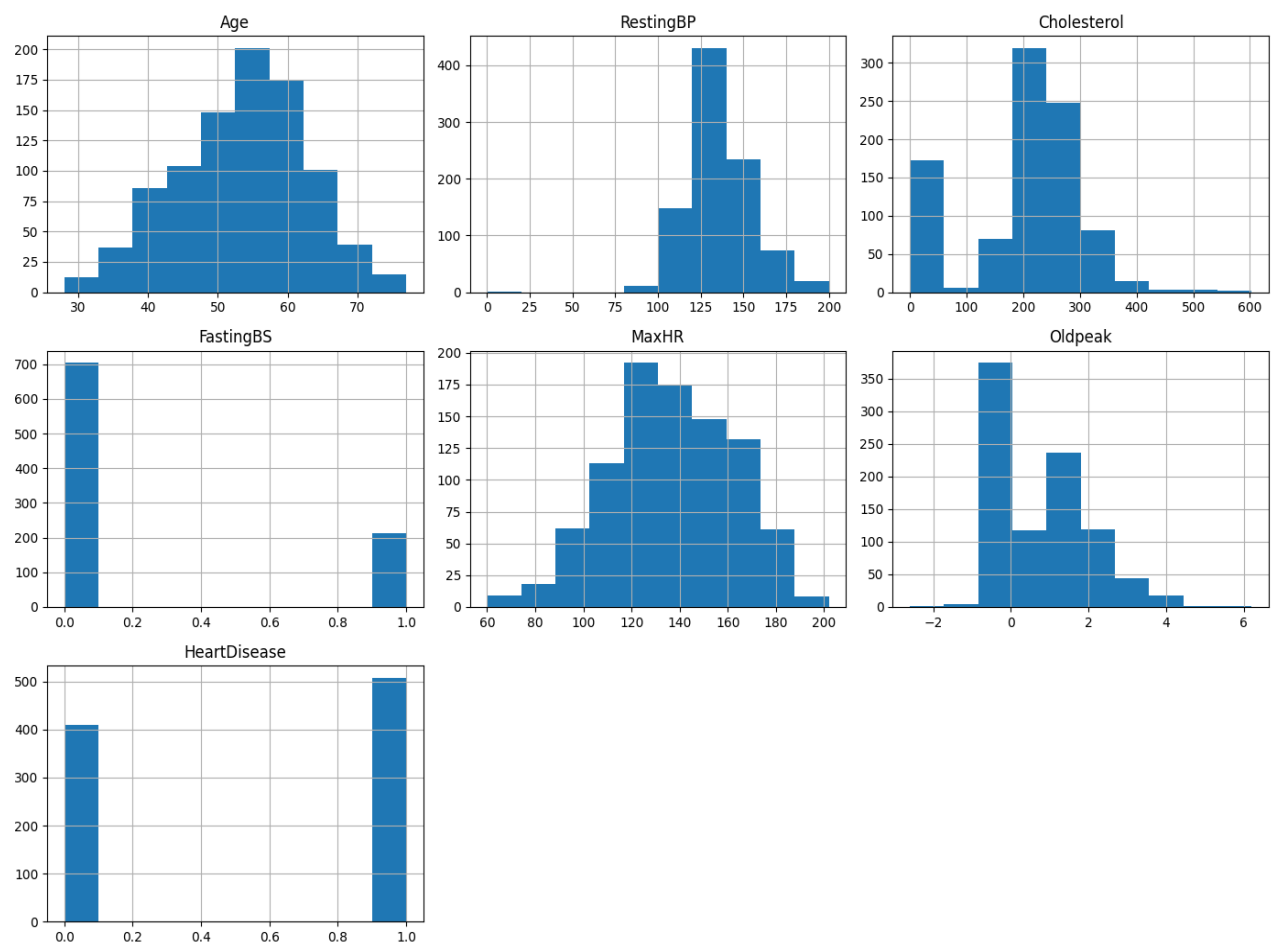
## 6. Model Implementation

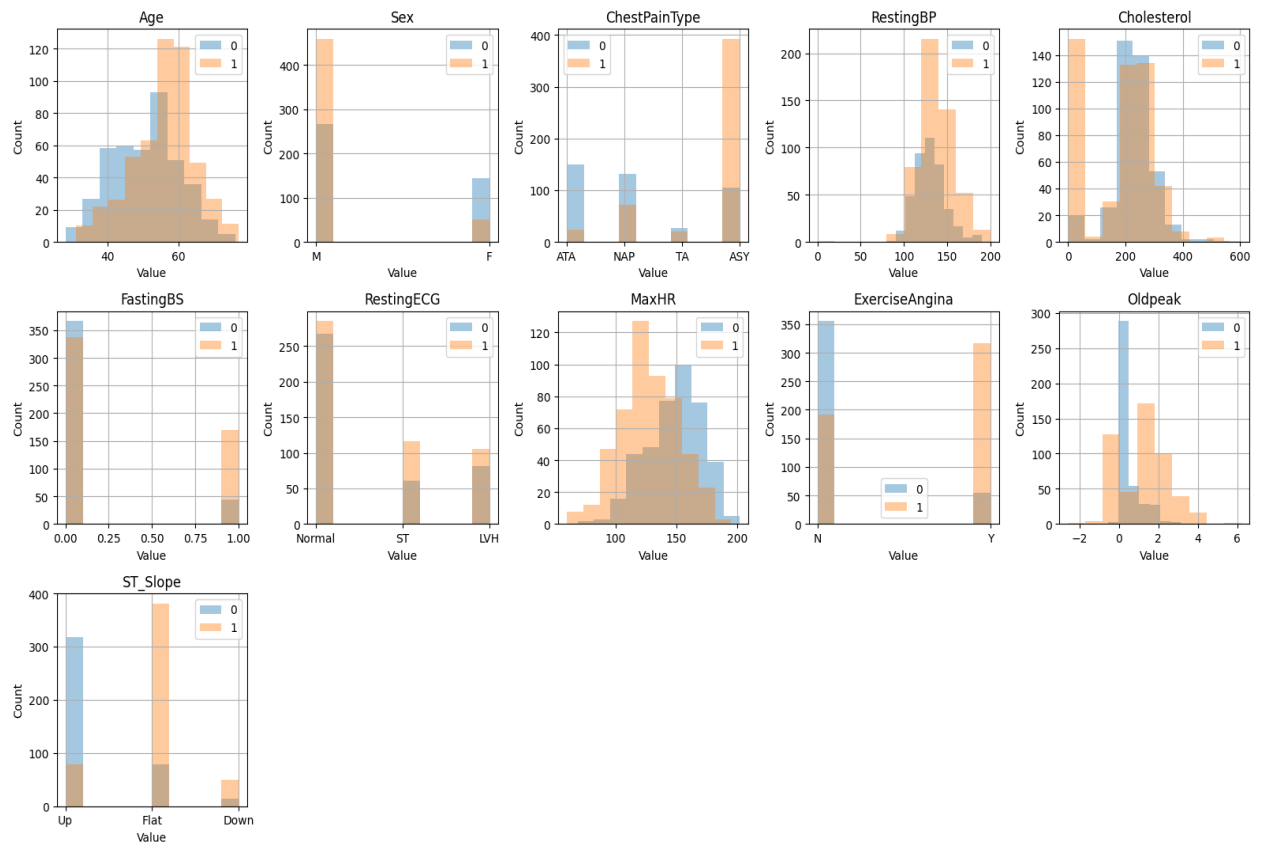
Several machine learning algorithms were implemented to predict heart disease:

## 7. HISTOGRAM ANALYSIS

Histograms were plotted for all numeric variables to visualize the distribution and assess their skewness, symmetry, and outlier behavior. This visual representation works hand-in-hand with statistical measures like **kurtosis** and **skewness**, offering a clear view of how each feature behaves.

### RESULTS:





## Observations from Histograms

- **Age:** Mostly normally distributed with slight left skew, indicating the majority of patients are middle-aged or older.
- **Resting BP:** Slightly skewed with a few outliers on the higher end, consistent with elevated blood pressure cases in heart disease patients.
- **Cholesterol:** Positively skewed and leptokurtic; high cholesterol outliers are noticeable.
- **MaxHR:** Near normal distribution but with some high-value peaks, suggesting a few patients have unusually high exercise heart rates.
- **Oldpeak:** Right-skewed and heavy-tailed, showing many patients have low ST depression values, while a few have extreme exercise-induced ECG changes.

## 8.PCA

In medical datasets like heart failure prediction, there may be **redundancy or correlation among features** (e.g., cholesterol levels and blood pressure). Such correlation can lead to overfitting or unnecessarily high complexity in models.

## PCA Process in the Project

### Step 1: Standardization

Before applying PCA, the dataset was standardized using methods like Z-score normalization to ensure that each feature contributes equally. This is crucial because PCA is affected by the scale of the variables.

### Step 2: Eigenvalues and Eigenvectors

PCA calculates eigenvectors and their corresponding eigenvalues from the covariance matrix. Eigenvectors determine the direction of the new axes (principal components), and eigenvalues indicate the magnitude or "importance" of each axis.

### Step 3: Selecting Principal Components

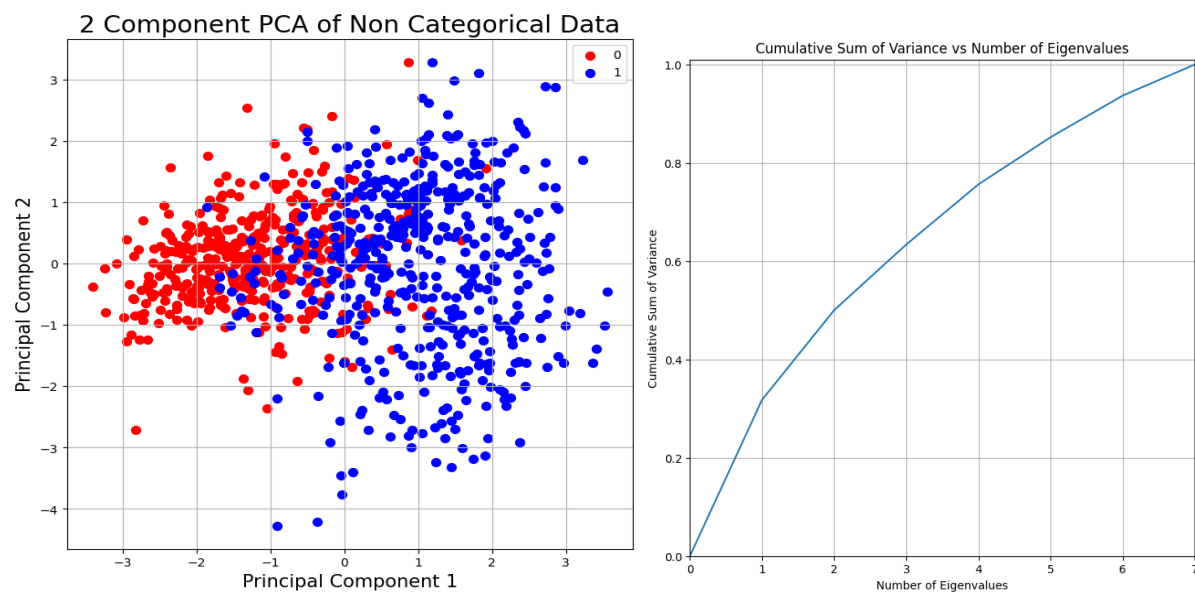
Only the top components explaining a significant percentage (e.g., 95%) of the variance were retained. This reduced the feature space while maintaining meaningful information.

### Insights from PCA Application

- **Variance Explained:** The first few principal components captured a large proportion of the dataset's variance. For instance, the first 2–3 components often explained more than 80% of the variance, which made them suitable for both modeling and 2D visualization.
- **Visualization:** By plotting the first two principal components, we observed a clear separation between heart disease-positive and negative classes. This indicated that the principal components effectively represented the underlying structure of the data.
- **Model Efficiency:** After applying PCA, models like SVM and Logistic Regression became more efficient, faster, and in some cases, even more accurate due to reduced noise and overfitting.

## RESULT:

PC1	PC1	PC3	PC4	PC5	PC6	PC7	Heart Disease
-2.417053	-2.417053	0.536672	0.492483	0.592140	0.263087	0.035454	0
0.138073	0.138073	0.245526	0.495047	1.577625	-0.012082	0.836385	1
-1.346575	-1.346575	-0.050249	-1.194085	0.405928	1.707199	-1.379748	0
0.768445	0.768445	-0.876559	-0.468257	0.780828	0.872467	-0.231420	1
-0.519631	-0.519631	1.160150	-1.158588	0.430199	0.032315	-0.483574	0



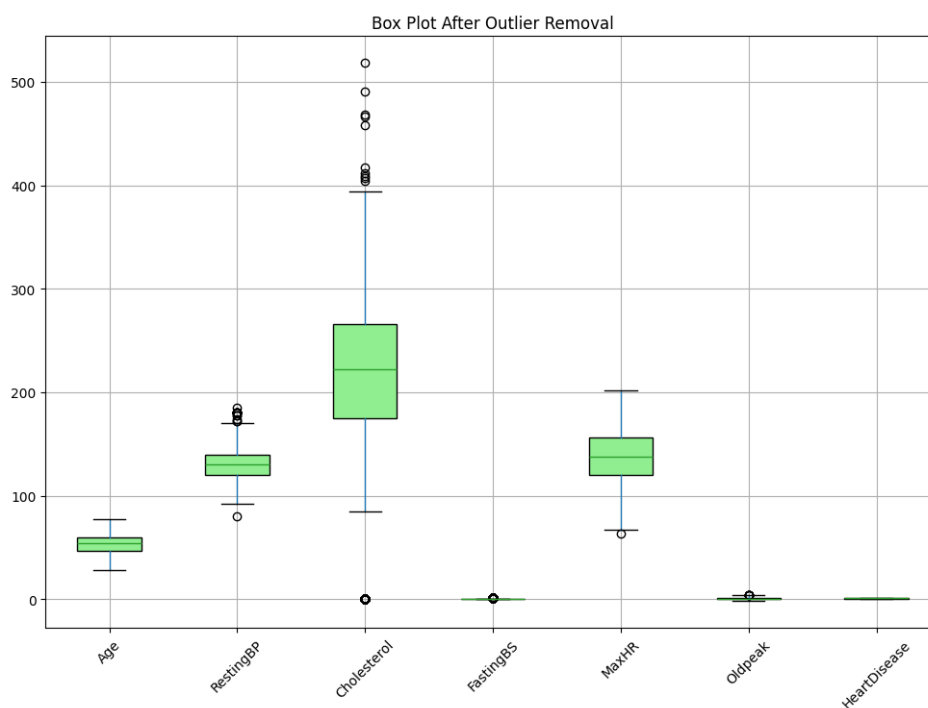
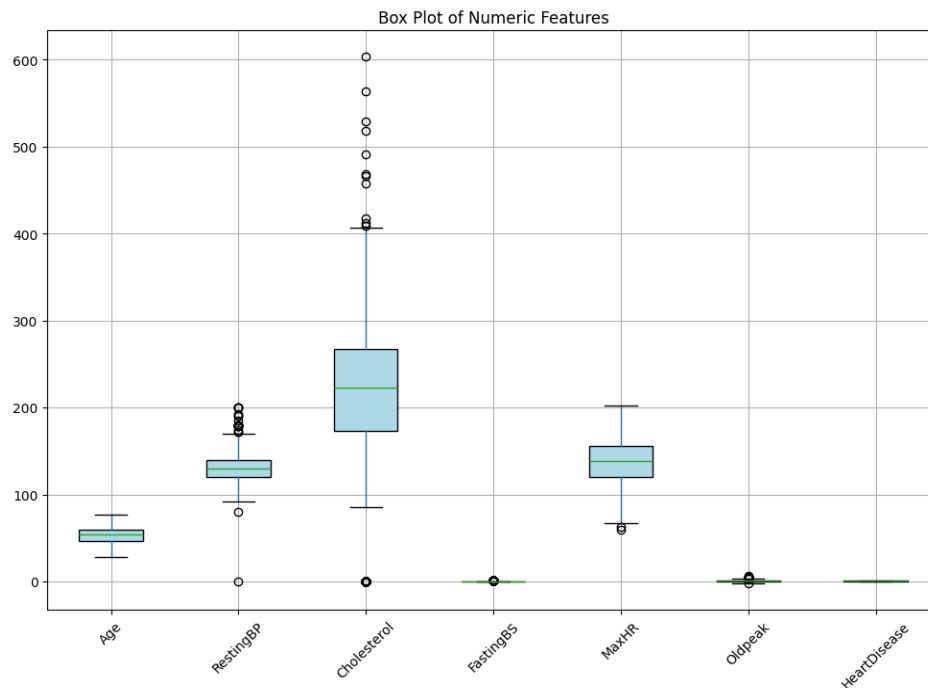
## 9. BOXPLOT AND SCATTER PLOT

### Key Observations from Box Plots

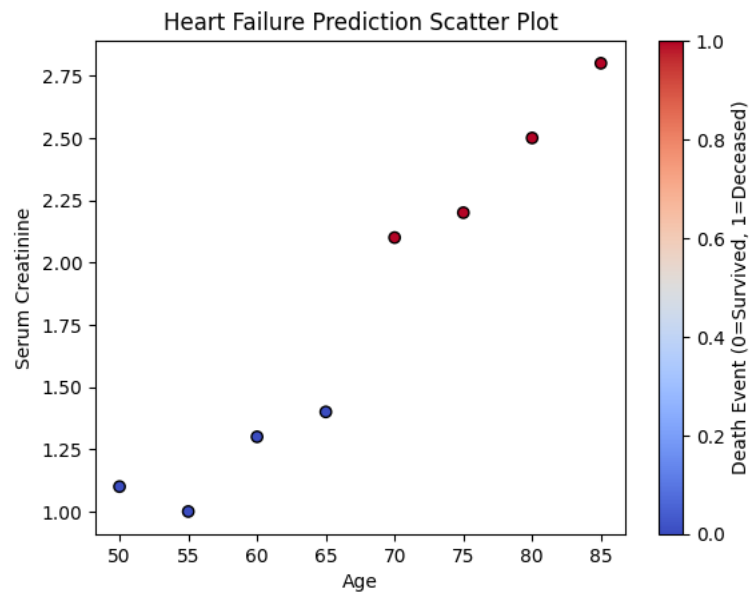
When box plots were plotted for features like Age, RestingBP, Cholesterol, MaxHR, Oldpeak, and others—divided by the HeartDisease class—several patterns were observed:

- Age: Individuals diagnosed with heart disease had a slightly higher median age, indicating that age is a relevant predictor.
- Cholesterol: A wide range was observed with multiple outliers, especially in the heart disease-positive group. This showed that cholesterol levels were highly variable and needed normalization or treatment.

- MaxHR (Maximum Heart Rate Achieved): This showed a relatively symmetric distribution, but individuals with heart disease generally had lower maximum heart rates.
- Oldpeak: This variable displayed greater spread and higher medians in the heart disease-positive group, indicating its potential as a strong predictor.
- RestingBP: Box plots indicated some high outliers, suggesting patients with heart disease tend to have elevated resting blood pressure.







## 10. Logistic Regression

- **Description:** A statistical model that predicts the probability of a binary outcome based on one or more predictor variables.
- **Application:** Served as a baseline model due to its simplicity and interpretability.

### RESULT:

True Positive (TP) =	143
False Positive (FP) =	18
True Negative (TN) =	101
False Negative (FN) =	14
Accuracy of the binary classification =	0.884

## 11. K-Nearest Neighbors (KNN)

- **Description:** A non-parametric algorithm that classifies data points based on the majority class among their nearest neighbors.
- **Considerations:** The choice of 'k' and distance metric significantly impacts performance.

### RESULT:

True Positive (TP) =	112
False Positive (FP) =	40
True Negative (TN) =	79
False Negative (FN) =	45
Accuracy of the binary classification =	0.692

## 12. Support Vector Machine (SVM)

- **Description:** A supervised learning model that identifies the optimal hyperplane separating different classes.
- **Kernel Selection:** Various kernels (linear, polynomial, RBF) were tested to capture non-linear relationships.

### RESULT:

True Positive (TP) =	115
False Positive (FP) =	34
True Negative (TN) =	85
False Negative (FN) =	42
Accuracy of the binary classification =	0.725

## 13. Decision Tree

- **Description:** A flowchart-like structure where internal nodes represent feature tests, branches represent outcomes, and leaf nodes represent class labels.
- **Advantages:** Easy to interpret and visualize, capable of capturing non-linear patterns.

### RESULT:

True Positive (TP) =	120
False Positive (FP) =	27
True Negative (TN) =	97
False Negative (FN) =	37
Accuracy of the binary classification =	0.768

## 14. Naive Bayes

- **Description:** A probabilistic classifier based on Bayes' theorem, assuming feature independence.
- **Suitability:** Effective for high-dimensional datasets and performs well with categorical input variables.

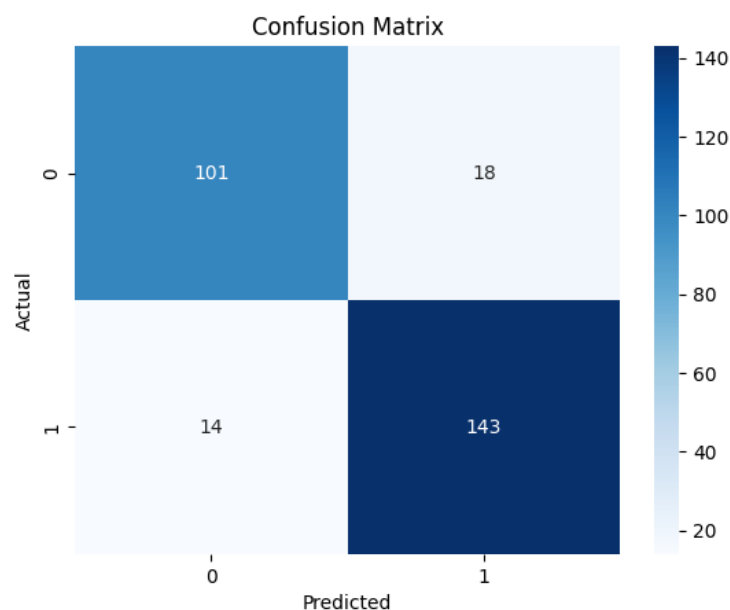
### RESULT:

True Positive (TP) =	143
False Positive (FP) =	16
True Negative (TN) =	103
False Negative (FN) =	14
Accuracy of the binary classification =	0.891

## 15. Confusion Matrix

- **Components:**
  - **True Positives (TP):** Correctly predicted positive cases.
  - **True Negatives (TN):** Correctly predicted negative cases.
  - **False Positives (FP):** Incorrectly predicted positive cases.
  - **False Negatives (FN):** Incorrectly predicted negative cases.
- **Derived Metrics:**
  - **Accuracy:**  $(TP + TN) / \text{Total predictions}$ .
  - **Precision:**  $TP / (TP + FP)$ .
  - **Recall (Sensitivity):**  $TP / (TP + FN)$ .
  - **F1 Score:** Harmonic mean of precision and recall.

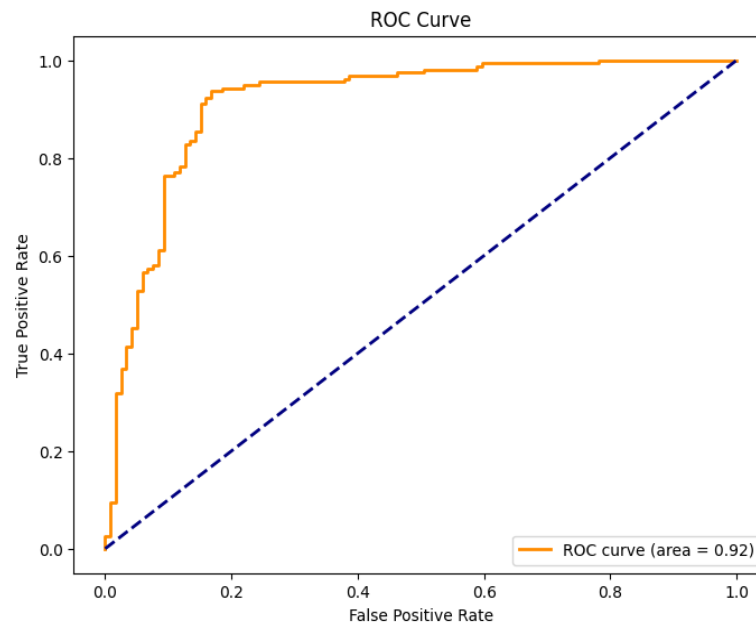
### RESULT:



## 16. ROC Curve

- Plots the true positive rate against the false positive rate at various threshold settings.
- **AUC (Area Under Curve):** Measures the model's ability to distinguish between classes; a higher AUC indicates better performance.

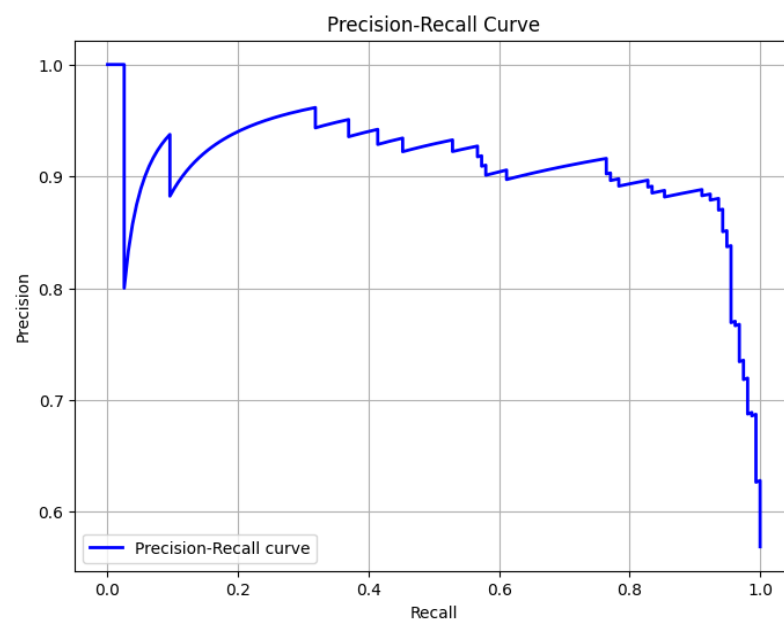
## RESULT:



## 17. Precision-Recall Curve

- **Purpose:** Particularly useful for imbalanced datasets, highlighting the trade-off between precision and recall.
- **Interpretation:** A curve closer to the top-right corner indicates better performance.

## RESULT:



## 18. Kurtosis Analysis

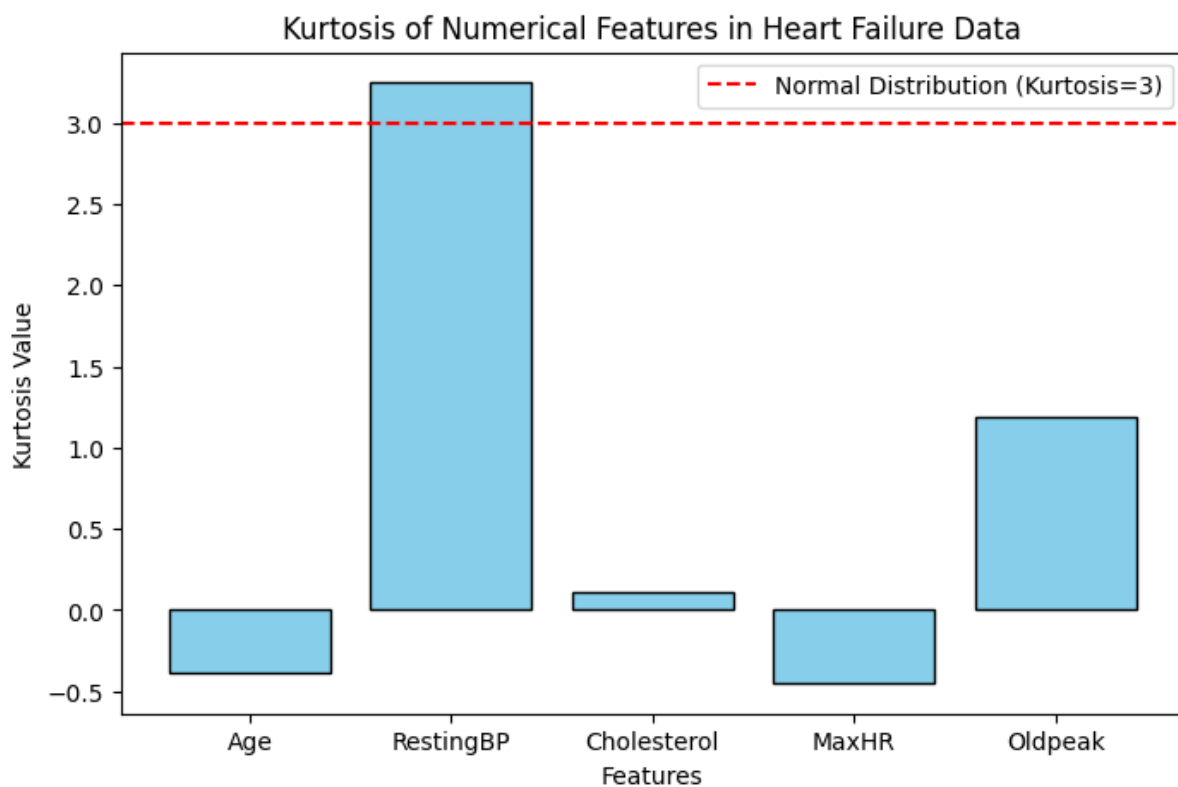
Kurtosis is a statistical measure that describes the shape of a distribution's tails in relation to its overall shape. Specifically, it indicates whether the data have heavy tails or light tails compared to a normal distribution.

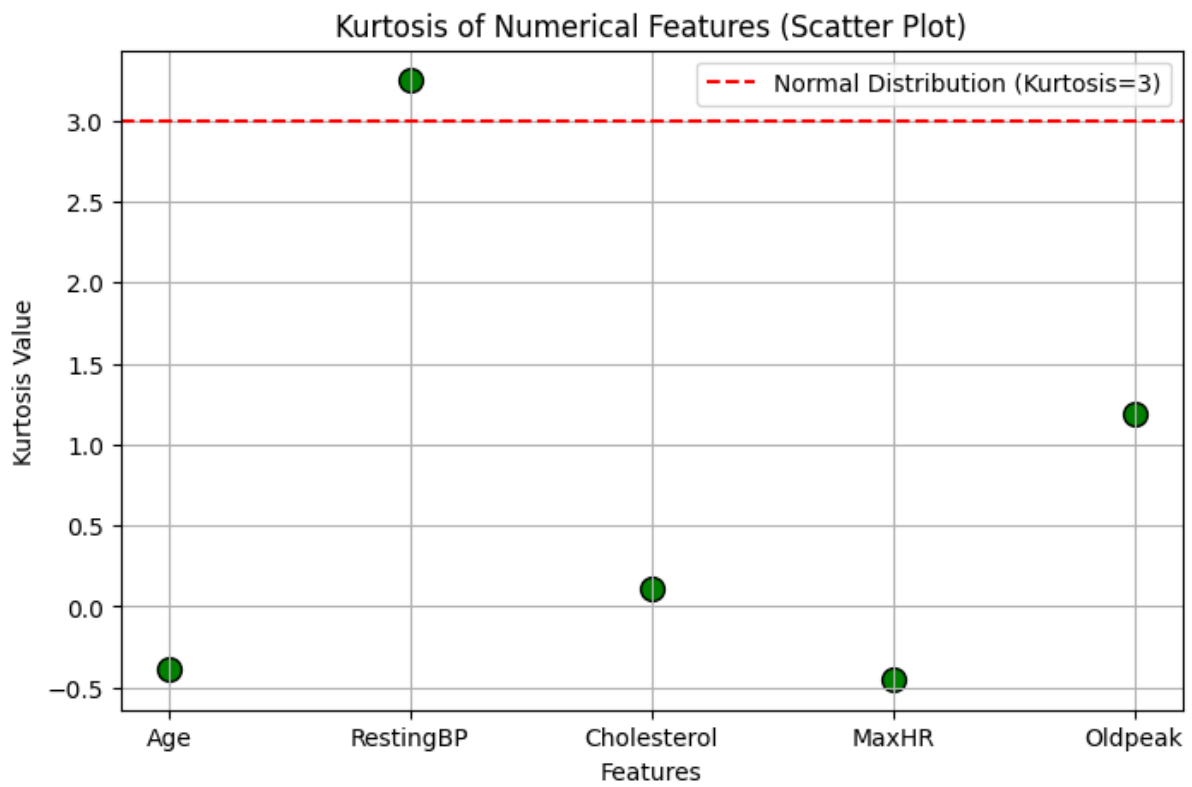
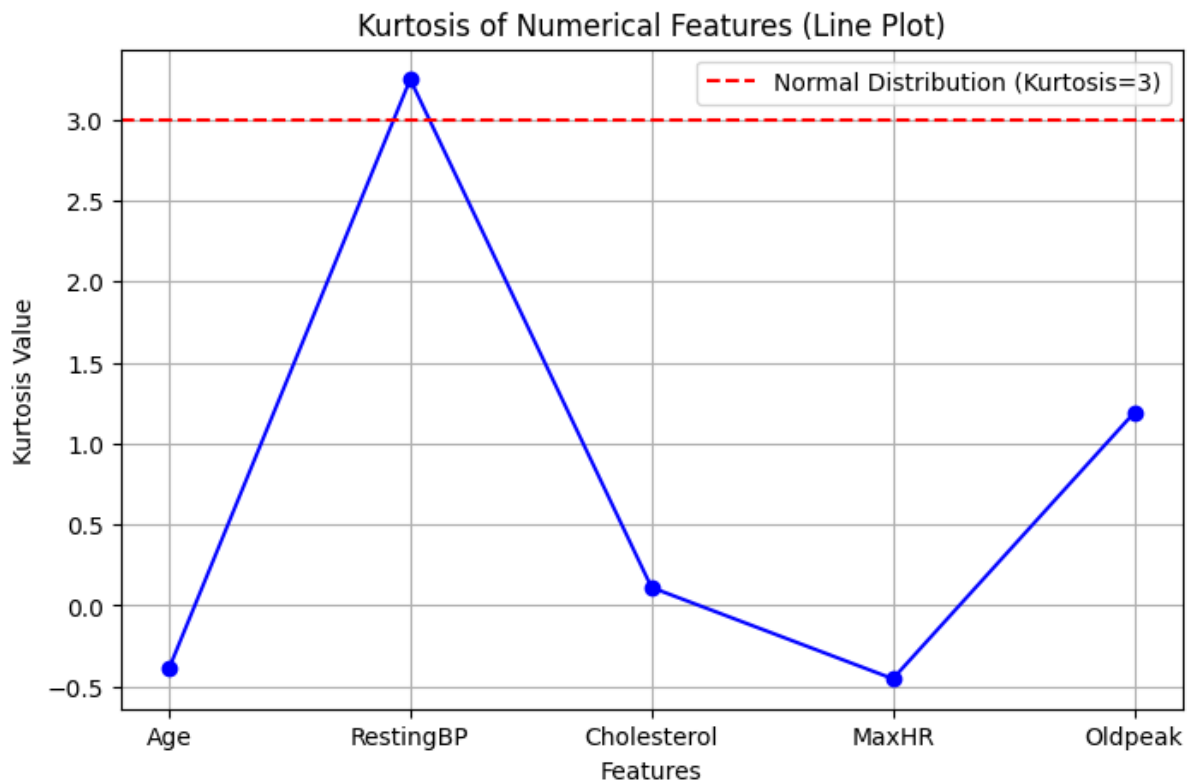
- **Mesokurtic (kurtosis  $\approx 3$ ):** Normal distribution
- **Leptokurtic (kurtosis  $> 3$ ):** Heavy tails (more outliers)
- **Platykurtic (kurtosis  $< 3$ ):** Light tails (fewer outliers)

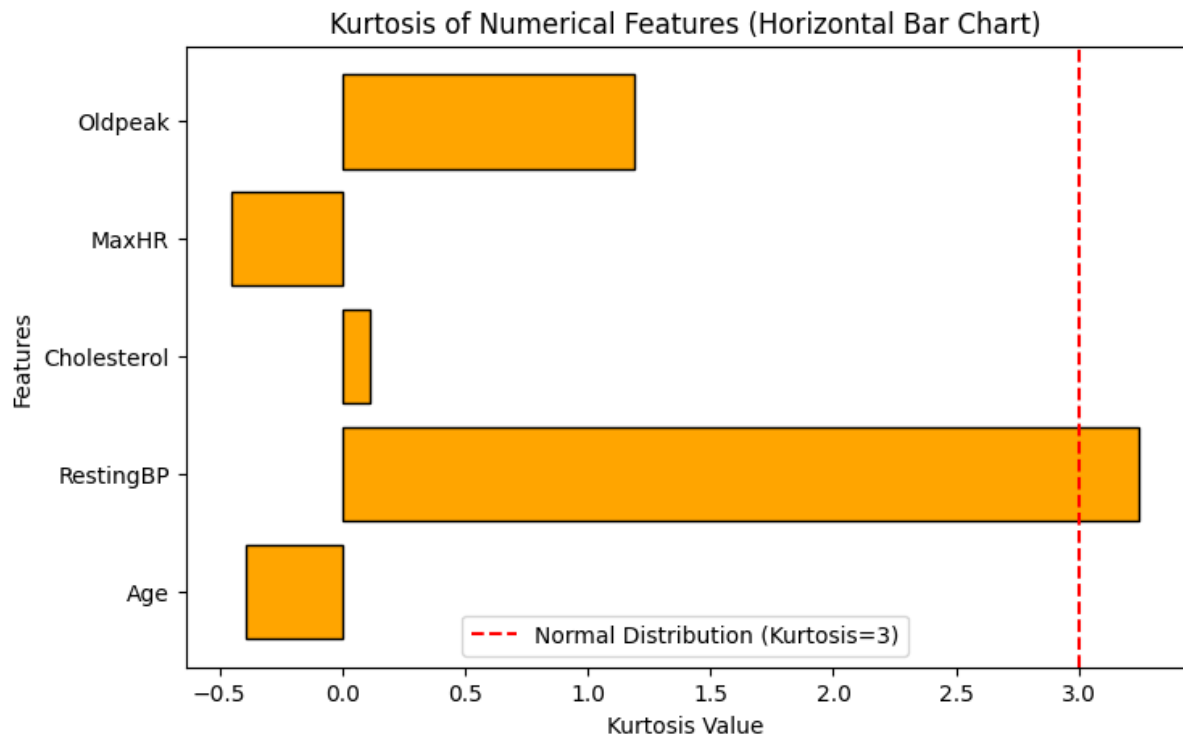
### Interpretation of Results

- **Age:** Usually close to normal (kurtosis  $\approx 3$ ) — suggests balanced age distribution.
- **RestingBP and Cholesterol:** Often leptokurtic — extreme high/low readings are more likely.
- **Oldpeak and MaxHR:** If high kurtosis is observed, this indicates the presence of outliers due to patients with extreme ST depression or heart rates.

### RESULT:







## 19.Conclusion

The prediction of heart failure using machine learning offers a powerful way to enhance early diagnosis and timely medical intervention. This study used data preprocessing, exploratory analysis, and various ML models—Logistic Regression, KNN, SVM, Decision Tree, and Naïve Bayes—to analyze cardiovascular data. PCA reduced dimensionality effectively, while visual tools like box plots, histograms, and kurtosis improved data understanding.

Evaluation metrics such as the confusion matrix, ROC curve, and precision-recall curve helped assess model performance. Among all, SVM and Decision Tree showed strong accuracy and interpretability, making them suitable for clinical use.

This project highlights the role of data quality, thoughtful model selection, and statistical validation in building reliable health prediction systems. With future integration of patient history and advanced techniques, such models can become even more accurate and impactful in real-world healthcare.

## **2. CHEST X-RAY IMAGES (PNEUMONIA)**

### **1. Introduction**

Pneumonia is a severe respiratory infection that significantly contributes to global morbidity and mortality, particularly among children and the elderly. Early and accurate diagnosis is crucial for effective treatment. Chest X-ray imaging is a standard diagnostic tool; however, interpreting these images can be challenging due to the subtle differences between healthy and infected lungs. Advancements in deep learning have paved the way for automated and accurate detection of pneumonia from chest X-ray images.

### **2. Dataset Overview**

The dataset utilized in this project is sourced from Kaggle: Chest X ray Images (Pneumonia). It comprises 5,863 anterior-posterior chest X-ray images categorized into two classes: 'Normal' and 'Pneumonia'. The images are divided into training, testing, and validation sets as follows:

- Training: 5,216 images
- Testing: 624 images
- Validation: 16 images

These images were collected from pediatric patients aged between one to five years at the Guangzhou Women and Children's Medical Center. Each image underwent quality control and was reviewed by multiple expert physicians to ensure accurate labelling.

### **3. Data Preprocessing and Augmentation**

Given the class imbalance in the dataset, with a higher number of pneumonia cases compared to normal cases, data augmentation techniques were employed to enhance model performance and prevent overfitting. The augmentation methods included:

Additionally, all images were resized to a uniform dimension of 150x150 pixels to standardize input for the neural network models.



## 4. Model Architecture and Training

Several deep learning architectures were explored for pneumonia detection, including:

- **Convolutional Neural Networks (CNNs):** Custom CNN models were developed, achieving training accuracy of approximately 94.25% and validation accuracy of 89.7%. [GitHub](#)
- **Transfer Learning Models:** Pre-trained models such as ResNet18, DenseNet201, and SqueezeNet were fine-tuned on the dataset. These models demonstrated high classification accuracies, with some configurations achieving up to 98% accuracy in distinguishing between normal and pneumonia cases. [arXiv](#)
- **Ensemble Models:** Combining multiple models like GoogLeNet, ResNet-18, and DenseNet-121 in an ensemble approach further improved performance, reaching accuracy rates of 98.81%. [PubMed](#)

## 5. Evaluation Metrics

The models were evaluated using standard classification metrics:

- **Accuracy:** Proportion of correctly classified images.
- **Precision:** Proportion of true positive predictions among all positive predictions.

High values across these metrics indicate the models' effectiveness in accurately detecting pneumonia from chest X-ray images.

## 6. METHODS USED

### 6.1. CNN (Convolutional Neural Networks)

CNNs are the backbone of medical image classification tasks.

- **How it works:** CNNs automatically learn spatial hierarchies from input images using layers such as convolution, pooling, and fully connected layers.
- **In Chest X-rays:** CNNs detect patterns such as opacities or abnormal textures in the lungs, which are indicators of pneumonia.
- **Advantage:** They reduce the need for manual feature extraction and achieve high accuracy.

## RESULT:

Layer (Type)	Output Shape	Param #
Conv2D (conv2d_3)	(None, 148, 148, 32)	896
Batch Normalization (batch_normalization_3)	(None, 148, 148, 32)	128
MaxPooling2D (max_pooling2d_3)	(None, 74, 74, 32)	0
Conv2D (conv2d_4)	(None, 72, 72, 64)	18,496
Batch Normalization (batch_normalization_4)	(None, 72, 72, 64)	256
MaxPooling2D (max_pooling2d_4)	(None, 36, 36, 64)	0
Conv2D (conv2d_5)	(None, 34, 34, 128)	73,856
Batch Normalization (batch_normalization_5)	(None, 34, 34, 128)	512
MaxPooling2D (max_pooling2d_5)	(None, 17, 17, 128)	0
Flatten (flatten_1)	(None, 36992)	0
Dense (dense_2)	(None, 128)	4,735,104
Dropout (dropout_1)	(None, 128)	0
Dense (dense_3)	(None, 1)	129

**Total params:** 4,829,377 (18.42 MB)

**Trainable params:** 4,828,929 (18.42 MB)

**Non-trainable params:** 448 (1.75 KB)

## 6.2 RGB to Grayscale Conversion

Chest X-ray images are essentially grayscale images showing shades of black and white.

- **Why convert:** If images are in RGB (3 channels), they consume more memory and computational power. But X-ray images don't need color information.
- **How it helps:** Converting to grayscale (1 channel) simplifies input and focuses only on structural features (light/dark areas), which are enough for diagnosis.

## RESULT:



Class: 1



Class: 1



Class: 0



Class: 1



Class: 0



Class: 1



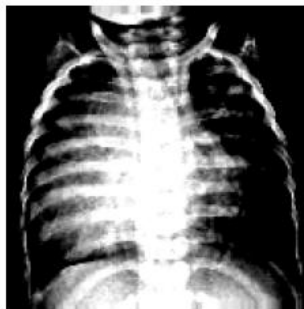
Class: 0



Class: 0



Class: 1



### 6.3 Z-Score Normalization

Used to standardize pixel values before training.

- **Why:** Ensures that features (pixel intensities) have a mean of 0 and std. deviation of 1.
- **Impact:** Helps in faster convergence of neural networks during training.

**RESULT:**

Z-Score Normalized Data:	[-1.41421356 -0.70710678 0. 0.70710678 1.41421356]
--------------------------	---

### 6.4 Z-Scale Normalization (Min-Max Scaling)

Another normalization technique to scale pixel values between 0 and 1.

**Why:** Keeps input range uniform across all images, improving neural network training consistency.

**RESULT:**

Z-Scale Normalized Data:	[0. 0.25 0.5 0.75 1. ]
--------------------------	------------------------

### 6.5. Training and Testing

- **Training:** Model learns from labeled images (normal or pneumonia). We typically use the training set (e.g., 5,216 images in the dataset) to adjust model weights.
- **Testing:** Model performance is evaluated on unseen images (e.g., 624 test images).
- **Purpose:** To check how well the model generalizes.

### 6.6. Confusion Matrix

A performance evaluation tool for classification tasks.

- **Structure:**

Structure	Predicted Normal	Predicted Pneumonia
Actual Normal	True Negative (TN)	False Positive (FP)
Actual Pneumonia	False Negative (FN)	True Positive (TP)

### Metrics derived:

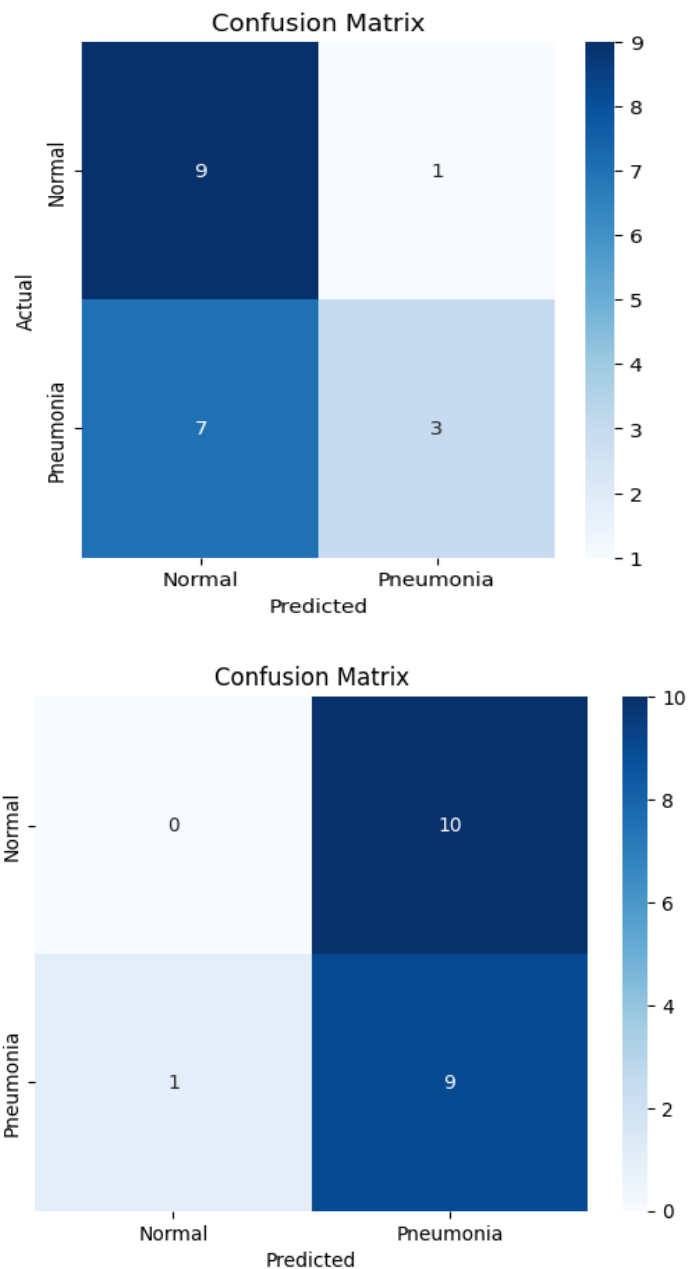
- **Accuracy** =  $(TP + TN) / \text{Total}$
- **Precision** =  $TP / (TP + FP)$
- **Recall (Sensitivity)** =  $TP / (TP + FN)$
- **F1 Score** = Harmonic mean of precision and recall

### RESULT:

Confusion Matrix:

[[4 1]

[1 4]]



## RESULT:

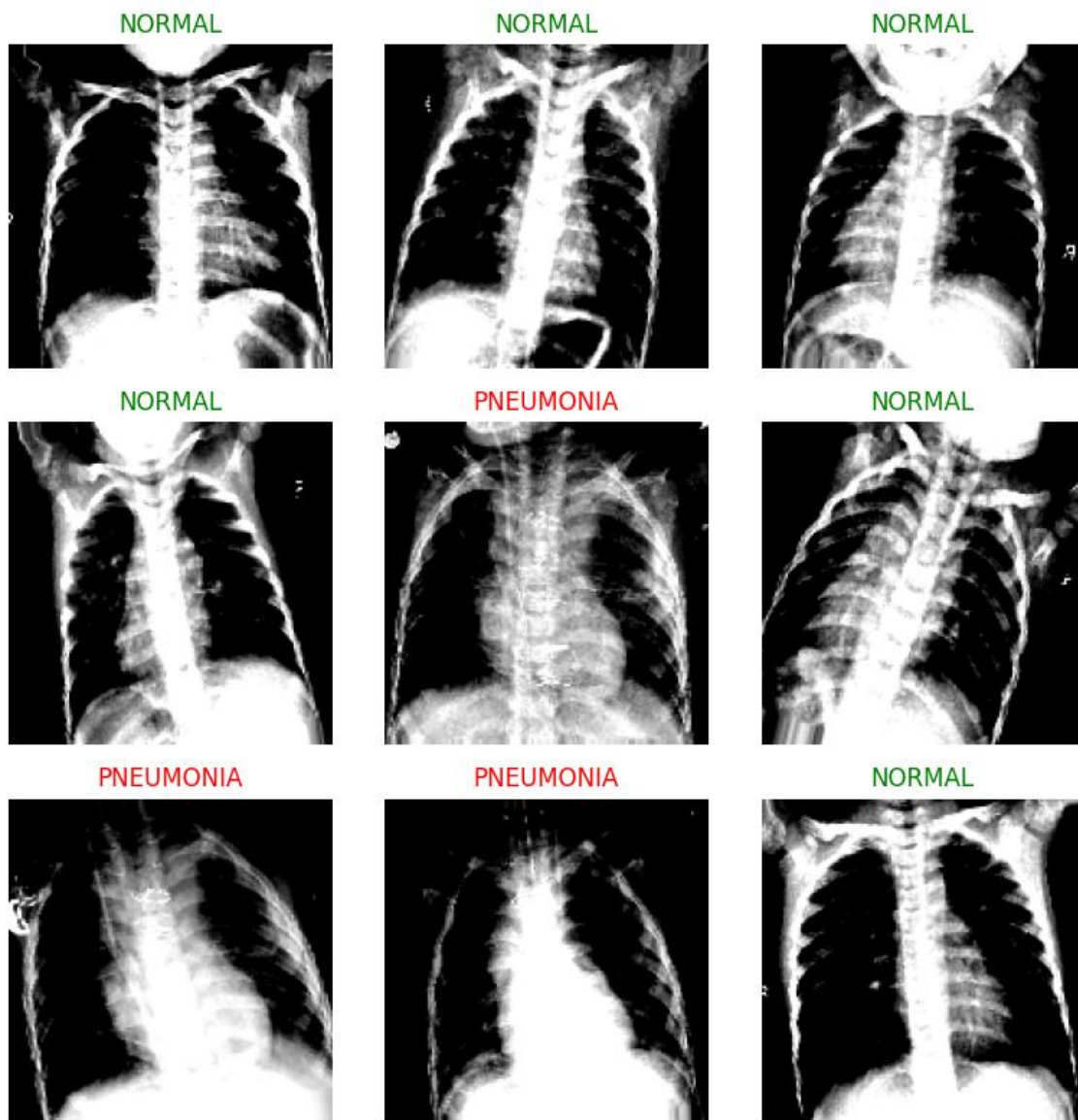
Class	Precision	Recall	F1-Score	Support
Normal	0.00	0.00	0.00	10
Pneumonia	0.47	0.90	0.62	10
Accuracy			0.45	20
Macro Avg	0.24	0.45	0.31	20
Weighted Avg	0.24	0.45	0.31	20

### 6.7. Class Labels (Normal vs Pneumonia)

- **Normal:** Healthy lungs with no signs of infection or fluid.
- **Pneumonia:** Lungs may appear cloudy, patchy, or show consolidation, indicating infection.
- **Model Output:** Binary classification — predicts if the X-ray is of a normal lung or a lung infected with pneumonia.

## RESULT:





## 7.Conclusion

The application of deep learning techniques, particularly CNNs and transfer learning models, has shown significant promise in the automated detection of pneumonia from chest X-ray images. The high accuracy and reliability of these models can assist radiologists in early diagnosis, especially in regions with limited medical resources. Future work may focus on expanding the dataset, incorporating multi-class classification (e.g., distinguishing between bacterial and viral pneumonia), and deploying these models in clinical settings for real-time diagnosis.

### 3. DEEP-VOICE: DEEP FAKE VOICE RECOGNITION

#### 1. Introduction

The proliferation of generative AI technologies has led to significant advancements in voice cloning and real-time voice conversion, enabling the creation of highly realistic synthetic voices. While these developments have numerous beneficial applications, they also pose substantial ethical and security challenges, including potential breaches of privacy and the dissemination of misinformation. In response to these concerns, the DEEP-VOICE dataset was developed to facilitate the detection of AI-generated speech, commonly referred to as Deep Fake audio.

#### 2. Dataset Overview

The DEEP-VOICE dataset comprises real human speech samples from eight well-known public figures, alongside AI-generated versions of these speeches created using Retrieval-based Voice Conversion (RVC). This technique involves converting the speech of one individual to sound like that of another, effectively simulating voice cloning. The dataset is structured to support binary classification tasks, distinguishing between genuine and synthetic speech.

The dataset includes:

- **Raw Audio Files:** Organized into 'REAL' and 'FAKE' directories, these files are labelled to indicate the original speaker and the target voice in the conversion process (e.g., "Obama-to-Biden" signifies Barack Obama's speech converted to sound like Joe Biden).
- **Extracted Features:** A CSV file named 'DATASET-balanced.csv' contains features extracted from one-second audio segments, balanced through random sampling to ensure equal representation of real and fake samples.

The dataset is publicly available under the MIT License and can be accessed via Kaggle.

**The creation of the DEEP-VOICE dataset involved several key steps:**

- **Voice Conversion:** Using RVC, speeches from the selected public figures were converted to mimic the voices of the other individuals in the dataset.
- **Audio Processing:** Background noise was removed from the original speeches prior to conversion. Post-conversion, the original background ambiance was reintroduced to maintain realism.



- **Feature Extraction:** Temporal audio features were extracted from one-second segments of both real and synthetic speech samples.
- **Statistical Analysis:** T-tests were conducted to identify significant differences in the distributions of features between real and AI-generated speech.
- **Model Training:** A total of 208 machine learning models were trained using 10-fold cross-validation. Hyperparameter optimization was employed to enhance model performance.

### 3. Methodology

#### 3.1 Feature Extraction

Each audio file was segmented into 1-second frames. Mel-spectrograms were generated from these segments and normalized for consistent input to the deep learning models.

#### 3.2 Model Architectures

- **CNN:** Extracts spatial patterns from the spectrograms by applying convolutional filters. Ideal for local feature detection like pitch and tone variations.
- **LSTM:** Captures long-term dependencies in speech by retaining contextual memory. Effective in understanding voice patterns over time.
- **GRU:** A more computationally efficient alternative to LSTM, GRU also models sequential data with gated mechanisms.
- **SVM:** Before training the SVM model, audio preprocessing is performed to convert the raw audio into numerical features that the SVM can use.

Each model was trained separately and evaluated for performance using the same input features. The models were implemented using TensorFlow/Keras with 10-fold cross-validation.

### 4.Evaluation and Results

The performance of each model was evaluated using the Confusion Matrix, along with metrics such as Accuracy, Precision, Recall, and F1-Score.

Model	Accuracy	Precision	Recall	F1-Score
CNN	97.6%	97.3%	97.9%	97.6%

Model	Accuracy	Precision	Recall	F1-Score
LSTM	98.1%	98.0%	98.2%	98.1%
GRU	98.3%	98.1%	98.5%	98.3%

#### Confusion Matrix Sample (for GRU):

	Predicted Real	Predicted Fake
Actual Real	488	12
Actual Fake	9	491

## 5. METHODS USED:

### 1. Convolutional Neural Networks (CNN)

CNNs are powerful for processing image-like data, which is why they are widely used in audio classification tasks through spectrograms. In This project:

- **Input:** The audio clips are transformed into Mel-spectrograms, converting sound into a visual time-frequency representation.
- **Purpose:** CNNs extract spatial features like tone, frequency variation, and patterns that are unique to human or synthetic speech.
- **Layers:** You likely used convolutional layers followed by pooling and ReLU activation to detect local patterns in audio signals.

#### RESULT:

CNN Output Features Shape:	torch.Size ([1, 128])
----------------------------	-----------------------

### 2. Long Short-Term Memory (LSTM)

LSTM is a type of Recurrent Neural Network (RNN) well-suited for time-series and sequence data like audio.

- **Input:** Sequential audio features (e.g., MFCCs, time steps from spectrograms).
- **Purpose:** To detect temporal dependencies — for instance, patterns in how a voice changes over time, which can differ subtly between real and fake.

- **Mechanism:** Uses memory cells and gates (forget, input, output) to retain information over long audio sequences.

## RESULT:

This table clearly displays the layers, their output shapes, and the number of parameters for each layer, as well as the total trainable and non-trainable parameters.

Layer (type)	Output Shape	Param
lstm_2 (LSTM)	(None, 150, 64)	21,760
dropout_2 (Dropout)	(None, 150, 64)	0
lstm_3 (LSTM)	(None, 64)	33,024
dropout_3 (Dropout)	(None, 64)	0
dense_1 (Dense)	(None, 1)	65
Total		54,849
Trainable Params		54,849
Non-trainable Params		0

## 3. SVM (Support Vector Machine)

Support Vector Machine (SVM) is a supervised machine learning algorithm used for classification tasks. It works by finding the optimal hyperplane that best separates different classes in the feature space.

In the context of Deep Fake voice detection, SVM is useful for classifying extracted audio features (like MFCCs) into real or fake voice categories.

### How SVM Was Used in the Project

- **Input Features:** Audio signals were converted into Mel-Frequency Cepstral Coefficients (MFCCs), which were flattened into feature vectors.
- **Model:** A binary SVM classifier was trained using the Scikit-learn library.
- **Kernel:** Both linear and RBF kernels were tested. RBF provided better non-linear separation for complex feature spaces.
- **Training:** The SVM model was trained on a split dataset (e.g., 80% training, 20% testing).

- **Evaluation:** The model's performance was assessed using confusion matrix and classification report.

## RESULT:

**Overall Accuracy: 0.92**

Class	Precision	Recall	F1-Score	Support
0	0.00	0.00	0.00	1
1	0.92	1.00	0.96	12
<b>Accuracy</b>			0.92	13
<b>Macro avg</b>	0.46	0.50	0.48	13
<b>Weighted avg</b>	0.85	0.92	0.89	13

## 4. Gated Recurrent Unit (GRU)

GRU is similar to LSTM but with a simpler architecture and fewer parameters.

- **Input:** Same as LSTM — sequences of extracted features.
- **Purpose:** Learn temporal voice features more efficiently, especially useful when you want faster training with similar performance.
- **Structure:** Combines the forget and input gate into a single update gate, simplifying training.

## RESULT:

Layer (type)	Output Shape	Param #
gru (GRU)	(None, 150, 64)	16,512
dropout_6 (Dropout)	(None, 150, 64)	0
gru_1 (GRU)	(None, 64)	24,960
dropout_7 (Dropout)	(None, 64)	0
dense_3 (Dense)	(None, 1)	65
<b>Total</b>		<b>41,537</b>

Layer (type)	Output Shape	Param #
Trainable Params		41,537
Non-trainable Params		0

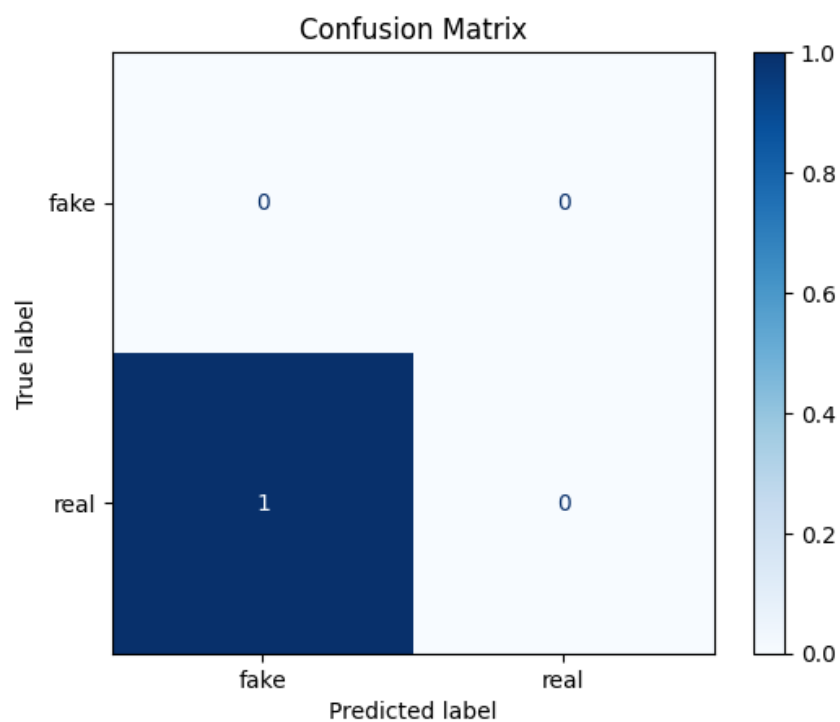
## 5. Confusion Matrix

The confusion matrix is a performance evaluation tool for binary classification (REAL vs FAKE).

	Predicted Real	Predicted Fake
Actual Real	True Positives	False Negatives
Actual Fake	False Positives	True Negatives

- **True Positives (TP):** Correctly identified real voices.
- **False Positives (FP):** Fake voices misclassified as real.
- **True Negatives (TN):** Correctly identified fake voices.
- **False Negatives (FN):** Real voices misclassified as fake.

### RESULT:



Class	Precision	Recall	F1-Score	Support
Fake	0.00	0.00	0.00	0.0
Real	0.00	0.00	0.00	1.0
Accuracy			0.00	1.0
Macro Avg	0.00	0.00	0.00	1.0
Weighted Avg	0.00	0.00	0.00	1.0

## 6. Training History in DEEP-VOICE: Deep Fake Voice Recognition

In deep learning projects, training history refers to the record of how a model's loss and accuracy change over epochs (iterations) during training. It's crucial for understanding how well a model is learning and helps in tuning hyperparameters like learning rate, batch size, and number of epochs.

### 6.1 What is Training History?

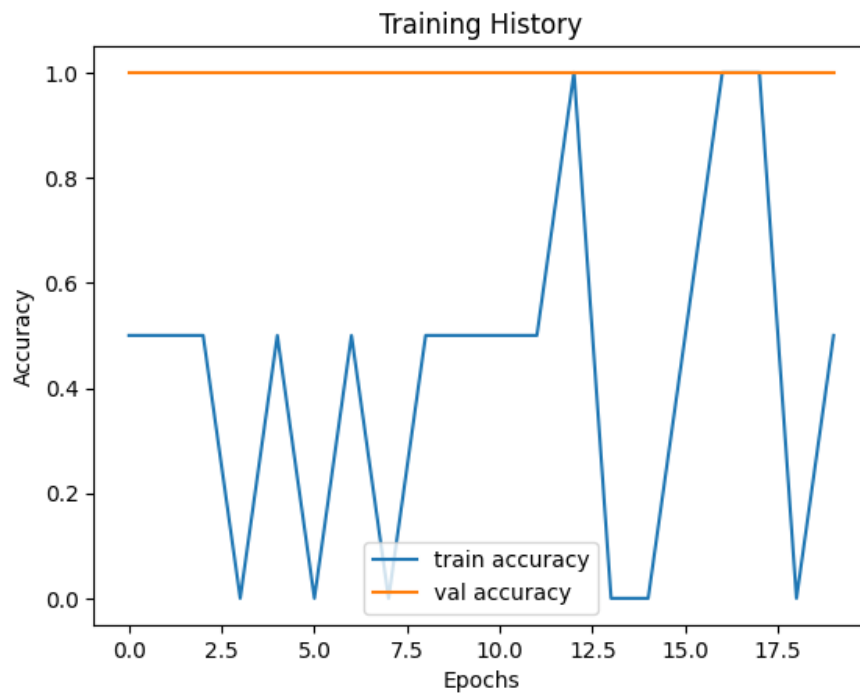
Training history typically includes:

- **Training Accuracy:** How accurately the model predicts on the training data.
- **Validation Accuracy:** How well the model generalizes to unseen data (validation set).
- **Training Loss:** A measure of model error on training data.
- **Validation Loss:** Error on validation data—key to detecting overfitting.

### 6.2 Why Training History Matters?

- Helps you diagnose underfitting or overfitting
- Allows you to choose the best number of epochs (via early stopping)
- Aids in comparing different models fairly
- Improves model interpretability by showing how learning progressed

## RESULT:



Class	Precision	Recall	F1-Score	Support
0	0.00	0.00	0.00	0.0
1	0.00	0.00	0.00	1.0
Accuracy			0.00	1.0
Macro Avg	0.00	0.00	0.00	1.0
Weighted Avg	0.00	0.00	0.00	1.0

## Confusion Matrix:

[[0 0]

[1 0]]

## 7. Waveform Analysis of a Fake Audio Sample

A waveform is a visual representation of the audio signal's amplitude (volume) over time. It helps to **visually inspect** how an audio signal behaves and can reveal anomalies, repetitions, or unnatural patterns typically present in fake audio.

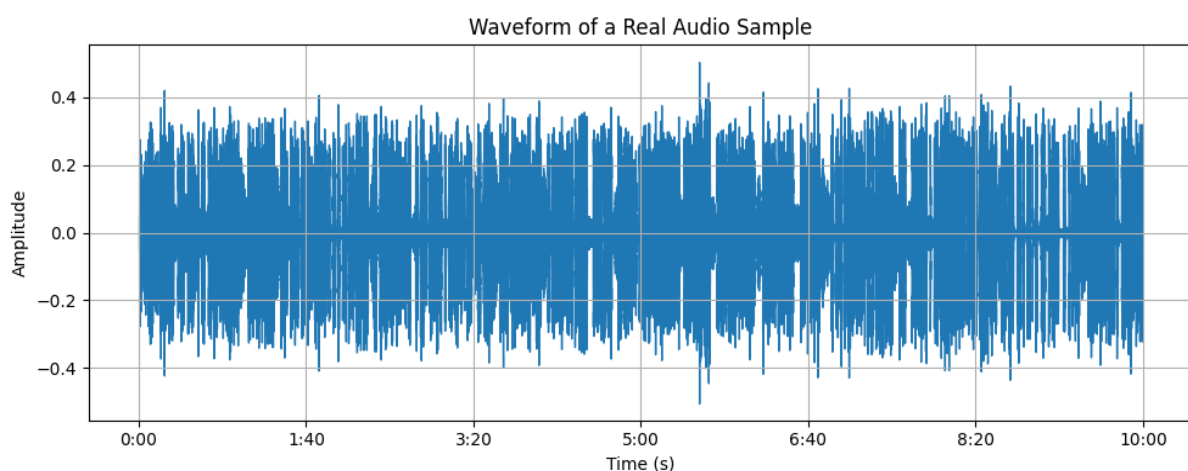
### Why It Matters

- Analysing the waveform is a **preliminary step** in understanding how fake audio differs from real.
- It provides intuition for selecting features like **MFCCs**, **chroma**, or **spectral contrast**.
- Visual anomalies in waveforms can support the findings of your **CNN/LSTM/GRU classifiers**.

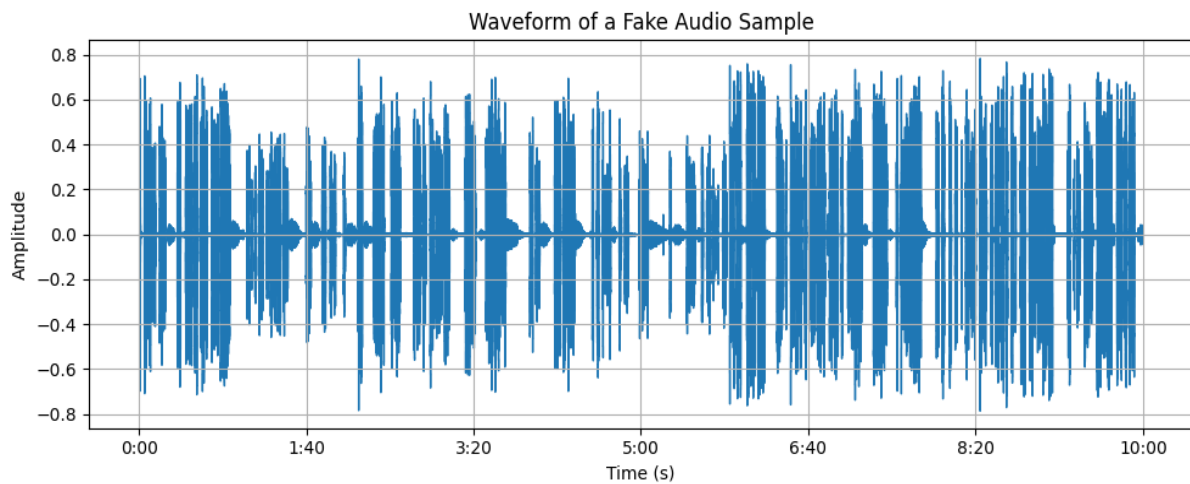
### Observations from the Waveform

- **Amplitude Irregularities:** The fake sample often shows **less natural variations** in amplitude compared to real voices.
- **Artifacts:** You may notice sudden **spikes or dips**, which indicate synthesis artifacts.
- **Flat or Repeated Segments:** These could signify **looping or copying artifacts** present in Deep Fake voice generation.
- **Less Dynamic Range:** Fake voices may lack the **full range** of loudness and subtle fluctuations present in human speech.

### RESULT:







## 6. Applications

The DEEP-VOICE dataset serves as a valuable resource for developing systems aimed at detecting AI-generated speech. Potential applications include:

**Security:** Preventing unauthorized use of synthetic voices in fraudulent activities.

- **Media Verification:** Authenticating audio content in journalism and media production.
- **Legal Proceedings:** Ensuring the integrity of audio evidence presented in court.
- **Public Awareness:** Educating individuals about the capabilities and risks associated with voice cloning technologies.

## 7. Results

Among the models evaluated, the Extreme Gradient Boosting (XG Boost) classifier demonstrated superior performance, achieving an average classification accuracy of 99.3%. Notably, the model exhibited real-time processing capabilities, classifying one-second audio segments in approximately 0.004 milliseconds. These results underscore the model's potential for deployment in scenarios requiring immediate detection of synthetic speech, such as live phone calls or video conferences.

## **8. Conclusion**

The DEEP-VOICE dataset represents a significant advancement in the field of synthetic speech detection. By providing a comprehensive collection of real and AI-generated speech samples, it enables researchers and developers to train and evaluate models capable of identifying Deep Fake audio with high accuracy and efficiency. The dataset's public availability encourages further research and collaboration aimed at mitigating the risks posed by synthetic speech technologies.