# Multi-Class and Binary Classification of Arrhythmias Using ECG Data from the MIT-BIH Arrhythmia and P-wave Databases

Jaladurgam Navya
Kent State Univeristy
jnavya@kent.edu

*Abstract*— This study explores the multi-classification of various types of arrhythmias using ECG data from MIT-BIH Arrhythmia Database, while another aim is to develop a model for binary classification based on a subset from MIT-BIH Arrhythmia Database, particularly P wave annotation dataset. The data set used for training the model involved P-wave annotation data only, prepared by experts who used SignalPlant. In this study, three different models used for MIT-BIH Arrhythmia Database are Random Forest, Support Vector Machine (SVM), and Neural Network. For MIT-BIH Arrhythmia Database, particularly P wave annotation dataset utilized four models are Random Forest, support vector machine, neural network, Gaussian Mixture Model (GMM). The study evaluates the model's performance by analyzing accuracy, recall and precision

*Keywords—segment, window, size, optimal, models*

## I. Introduction

Recently, the treatment and cure of all the ailments related to the heart has been a major field of research. Timely identification of heart problems is important especially for patients in the risk zone as the public health risk of sudden death due to heart problems is still real. Diagnosing critically ill ECG patients if they are referred to clinics is not limited to the regular spot checks.

The MIT-BIH Arrhythmia Database is essential to this type of investigation. This database is massive, with a wide variety of heart records and their diagnoses terminology aimed at enabling the researchers in coming up with algorithms for arrhythmia identification, developing and testing those algorithms .These are especially beneficial when improving the rate of diagnosis of interventional procedures, doing away with misconceptions. The MIT-BIH Arrhythmia Database P-Wave Annotations contribute to public health initiatives aimed at reducing sudden cardiac death and improving overall heart health through early detection and intervention. In this paper some of the methodologies are used to find the Arrhythmia types and binary classification of p-wave or not p-wave. Got 90% accuracy for random forest model for first database and 100% for second database.

## II. Database

### A. MIT-BIH Arrhythmia Database

For the MIT-BIH Arrhythmia Database the data is collected from the joint effort of Beth Israel Hospital and MIT from 1975 to 1979 resulted in the creation of the MIT-BIH Arrhythmia Database that has become a valuable addition for researchers investigating heart rhythmic disorders. Each of these 48 ECG recordings are on an average 30 minutes long and are available for analysis to 47 different subjects across multiple studies. These recordings feature an inpatients and outpatients includes a few outpatients among them. Twenty-Three out of the 48 recordings were randomly sampled. The database is made up of three types of files that is .dat files which contain ECG signals, .atr files which are detailed annotations of the heartbeats, .hea files which consist of information of the recordings and .xws file serves as a supplementary file that provides metadata about the ECG records.

### B. MIT-BIH Arrhythmia Database P-Wave Annotations

The MIT-BIH Arrhythmia Database with P-Wave Annotations the second database. This database has reference P-wave annotations for twelve signals from the MIT-BIH Arrhythmia Database. The signals were selected because they have conditions that make it harder to detect P-waves. This has .hea contain header information all separately, the .dat contain the real ECG signal while the .atr contains markings with p and null.

## III. Datapreprocessing

After loading three types of files. The header file has number of signals, sampling frequency, no of samples, signal names, adc gain, baseline, units. Next, the data preprocessing starts .

### A. MIT-BIH Arrhythmia Database

Replacing any missing values (NaN) in the signal with zeros and prepares to convert ADC (Analog-to-Digital Converter) values into physical units. This step involves adjusting the signal based on baseline and ADC gain values from the header. Scaling and normalization is performed. The normalized signal is then created by scaling the standardized data to a range between 0 and 1. The record number is printed as part of the preprocessing step when the .xws data is available. After preprocessing, the function detects the peaks and valleys in the signal .
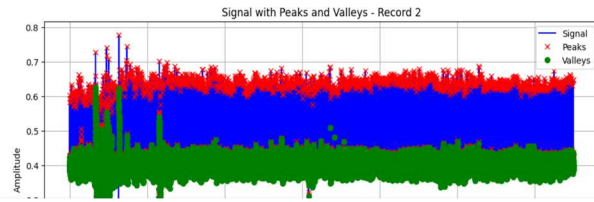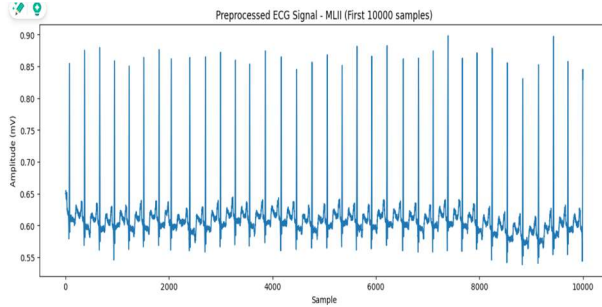
Fig 1: peaks and valleys



Fig 2: Preprocessed signal

The figure above is the preprocessed signal .The plot shows the first 10,000 samples of the ECG signal with the signal name (MLII) from the first record.

### B. MIT-BIH Arrhythmia Database P-Wave Annotations

After loading the data of 12 ECG records, headers and P-wave annotations. The data preprocessing steps are bandpass filter removes noise from the signal while keeping only the important frequencies between 0.5 Hz and 50 Hz. The filtered signal is then adjusted to have a mean of 0 and a standard deviation of 1. The code identifies the highest points (peaks) and lowest points (valleys) in the adjusted signal, ensuring that these points are spaced at least 150 samples apart.

```
Record 100:
    Original Signal Length: 1300000
    Filtered Signal Length: 1300000
    Scaled Signal Length: 1300000
    Number of Peaks Detected: 6168
    Number of Valleys Detected: 6771

Record 101:
    Original Signal Length: 1300000
    Filtered Signal Length: 1300000
    Scaled Signal Length: 1300000
    Number of Peaks Detected: 6136
    Number of Valleys Detected: 6928

Record 103:
    Original Signal Length: 1300000
    Filtered Signal Length: 1300000
    Scaled Signal Length: 1300000
    Number of Peaks Detected: 6386
    Number of Valleys Detected: 6253
```

Fig 3: Processed signal values

The before figure displays the some of the records output of preprocessing are original signal length, filtered signal length, scaled signal length, number of peaks detected and number of valleys detected.

Each record is made up of a fixed number of samples 1,300,000 samples in total. In this case, the data are uniformly captured within different records. The number of peaks ranges from a minimum of 5,573 to a maximum of 6,890, indicating differences in signal characteristics among other records. The number of valleys detected also shows considerable variation, from 6,091 to 7,347. Similar to peaks this variability may reflect different physiological conditions.
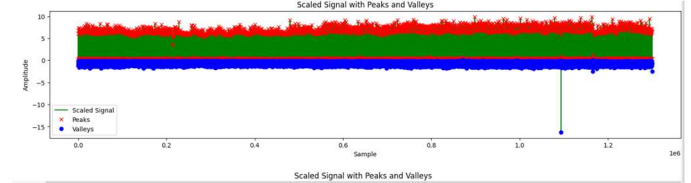


Fig 4: Scaled signal with peaks and valleys

## IV. SEGMENTATION

### A. MIT-BIH Arrhythmia Database

After preprocessing the next step is the segmentation. The steps in segmentation include segmenting the ECG signals, finding the optimal window size, and applying a sliding window technique. Segmentation for ECG signals, breaks an ECG signal up into parts at the desired window and step sizes which will overlap, and that function moves through the signal to extract each part and then move the next step length to start the process all over again. After the segmentation is done, the code uses this function in ECG signals that have been pre-processed multiple times but particularly that aim is sought for the first channel in the presence of more than one channel.

Next, find the optimal sized window to do the best segment of the ECG signals based on the standard deviation of the segments. the ECG signal with the minimum and maximum window sizes along with a step size. It looks through the set of potential window sizes, calculates the mean standard deviation of the segments imposed by the above method. If a better score is achieved, the 'best score' and its corresponding size will be updated. Sliding window technique using the optimal window size identified earlier for segmenting the ECG signals

```
Record 1: Optimal window size = 500
Record 2: Optimal window size = 500
Record 3: Optimal window size = 300
Record 4: Optimal window size = 500
Record 5: Optimal window size = 500
Record 6: Optimal window size = 500
Record 7: Optimal window size = 500
Record 8: Optimal window size = 500
Record 9: Optimal window size = 500
Record 10: Optimal window size = 500
Record 11: Optimal window size = 500
Record 12: Optimal window size = 500
Record 13: Optimal window size = 400
Record 14: Optimal window size = 500
Record 15: Optimal window size = 500
Record 16: Optimal window size = 500
Record 17: Optimal window size = 450
Record 18: Optimal window size = 500
Record 19: Optimal window size = 500
Record 20: Optimal window size = 500
Record 21: Optimal window size = 500
Record 22: Optimal window size = 450
Record 23: Optimal window size = 450
Record 24: Optimal window size = 500
Record 25: Optimal window size = 500
Record 26: Optimal window size = 500
Record 27: Optimal window size = 500
Record 28: Optimal window size = 500
Record 29: Optimal window size = 500
Record 30: Optimal window size = 500
Record 31: Optimal window size = 500
Record 32: Optimal window size = 500
Record 33: Optimal window size = 500
Record 34: Optimal window size = 500
```

Fig 5: Optimal window size

## B. MIT-BIH Arrhythmia Database P-Wave Annotations

The steps in segmentation include finding the optimal window size, creating sliding windows, optimal window segmentation and applying a sliding window technique. The purpose of these steps is to operate ECG signal by segmenting them to specific windows with respect to P-wave annotations and detecting the occurrence P-waves in them. The search for optimal window size focuses on the task of finding the most appropriate size of the ECG signal window based on the attributes of the P-wave annotations. It checks for the presence of at least two P-wave annotations to carry out this computation. The distinctions between consecutively occurring P-waves are defined. A central tendency and the variation of the distribution of these values helps to settle on the perfect window size; that is inclusive of specific limits between a lower end and upper end.

The step size is computed by sliding window by the specified overlap ratio. Each of the resulting segments' lengths shall not exceed the window width specified. The subsequent indices and events are stored together for further analysis.

Optimal window segmentation segments the ECG signal into optimal windows and counts the P-waves within each segment. For each window, it counts the number of P-wave annotations that fall within the window's boundaries and stores these counts. After processing all records, it finds the median of the optimal window sizes and prints the overall optimal window size in both samples and seconds. Overall, the optimal window size for the signal is based on the characteristics of the P-wave annotations. After determining the optimal window size, the signal is segmented, and the number of segments created. Calculating one optimal window size for all ECG records gives a single, consistent window length. This makes it easier to analyze each record in the same way, saves time, and reduces the work needed, especially when dealing with many records or limited computing power. In the figure below shows the overall optimal window size

Overall optimal window size: 797 samples (2.21 seconds)

Fig 6: Overall optimal window size

## V. FEATURE EXTRACTION

### A. MIT-BIH Arrhythmia Database

Performing feature extraction and analysis on the ECG signal , which involves gathering the relevant information, be it in statistical form. The identification of local maxima and minima, as well as computation of the amplitudes and the durations of the wave ensure that the understanding of the signal's behavior and can be used for further analysis. In feature extraction for each sample in the annotations, a segment of the signal defined by a window size of 500 . Various statistics are calculated for each of these segments, average, standard deviation, maximum, minimum , and range (max-min value).

The amplitudes and intervals find how high each peak is compared to the nearest low point and measure the distance between each peak. To determine the height of a peak, the value at the low point is subtracted from the value at the peak.
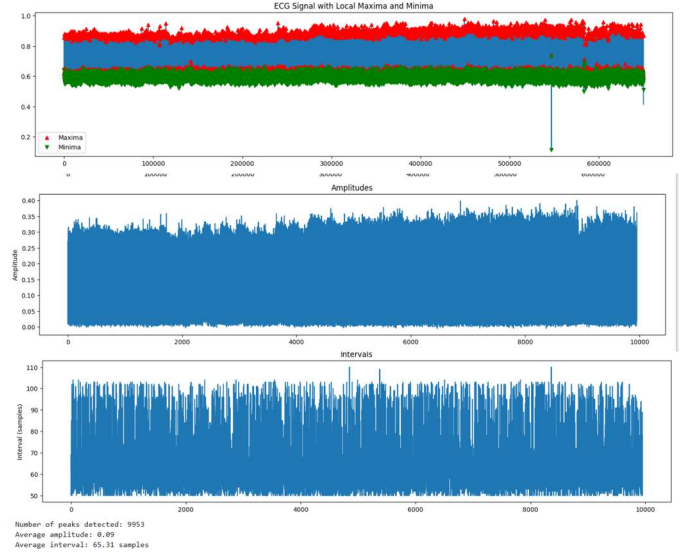


Fig:7 Local maxima and minima of first record

The above figure display's the number of detected peaks, the average amplitude, and the average interval for the first record. The number of peaks detected are 9953, the average amplitude is 0.09 and average interval is 65.31 samples.

### B. MIT-BIH Arrhythmia Database P-Wave Annotations

Finding local maxima and minima, calculate intervals between these points, and computing amplitude statistics. For each ECG record, it extracts the signal, identifies extrema, and calculates RR intervals and peak-to-peak amplitudes.
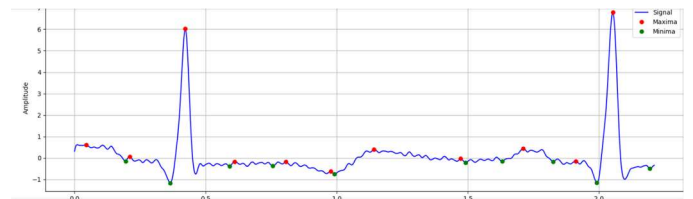


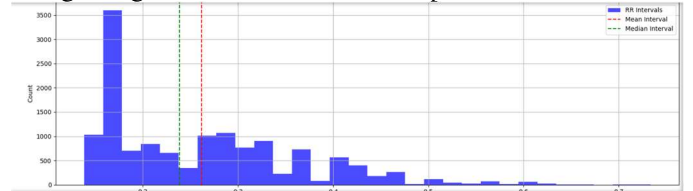Fig :7 Signal with detected Extrema points for record 100



Fig 8:Distribution of RR Intervals for record 100

Extracting features to analyze ECG signals by extracting time and frequency domain features. It calculates statistics like

mean and standard deviation, along with RR interval and amplitude features in the time domain and uses fast fourier transform to compute power in various frequency bands. In the below figure displays the extracted features.

```
Feature extraction completed!
        mean   std  skewness   kurtosis  rms  rr_mean   rr_std   rr_ratio \
0  8.449993e-18  1.0  4.847989  29.774128  1.0  0.261933  0.102690    5.28
1 -6.367556e-18  1.0  4.759590  28.476363  1.0  0.224251  0.079909    4.88
2  7.993606e-18  1.0  3.094536  17.463435  1.0  0.400462  0.116561    9.14
3 -4.110216e-18  1.0  2.334093  15.222323  1.0  0.245132  0.096285    6.66
4 -5.870176e-18  1.0  2.035554   8.569917  1.0  0.255824  0.092333    4.80

   rr_pnn50  peak_to_peak  peak_amplitude_mean  peak_amplitude_std \
0  64.678999     26.165330             1.233167            2.605211
1  53.440994     15.577545             0.941935            2.403072
2  79.913487     17.738517             1.674374            2.448331
3  56.782077     18.520535             0.846848            1.789099
4  64.458301     12.266346             0.790916            1.808874

      vlf_power     lf_power     hf_power  lf_hf_ratio   total_power
0  4445.344184  17449.340112  2.631129e+05    0.066319  2.850076e+05
1  14987.266126  45188.936690  5.117507e+05    0.088303  5.719269e+05
2  16457.666871  72116.739021  1.176809e+06    0.061282  1.265384e+06
3  26112.021396  98506.311547  8.611950e+05    0.114383  9.858134e+05
4  18888.611691  57410.757106  3.627597e+05    0.158261  4.390591e+05
Extracted features shape: (12, 17)
```
Fig 9: Feature Extraction

## VI. DATA VISUALIZATION FOR ARRHYTHMIA

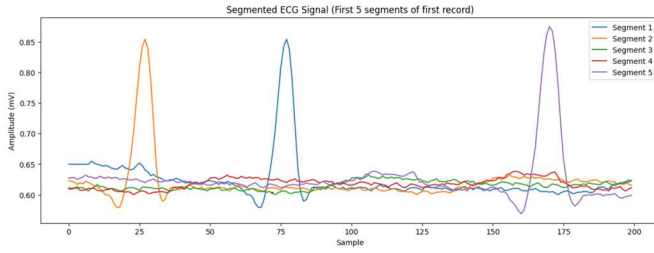### A. Segemented ECG signal for MIT-BIH Arrhythmia Database


Fig 10: Segmented ECG signal

The above figure presents the ECG signals extracted from the MIT-BIH Arrhythmia Database in a plot, with five parts (segments) from one ECG recording. Each segment of the five has a different color line in the figure. The x-axis represents the sample point which may signify time or the sampling rate; the y-axis shows the amplitude of the signal in millivolts (mV) revealing the electric activity of the heart. The huge positive peaks in the segments 1, 2, and 5 are the complex QRS and are indicative of the contraction of the ventricles of the heart. These QRS complexes are quite essential markers in ECG analysis since its magnitude, duration, and waveform shape are used to indicate the different pathological conditions of the heart.

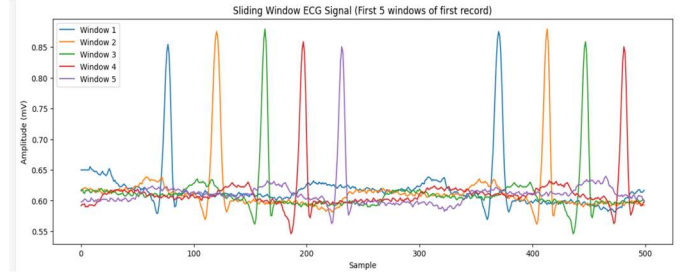### B. Sliding window ECG signal forMIT-BIH Arrhythmia Database


Fig: 11 Sliding window ECG signal for first 5 windows of first record

Partitions the ECG signal into smaller segments, each having a window of width called optimal size. The most important characteristics are the very apparent electrocardiogram graphs QRS (in this range sharp peaks in 0.85-0.87 mV can be seen and the level of the signal without the sharp peaks remains still around 0.60-0.62 mV in most cases). The recording captures two complete cardiac cycles within the sampling period, with each window showing the same fundamental pattern but shifted in time relative to the others. The regular spacing between peaks and consistent amplitude suggests a normal sinus rhythm, likely from a healthy heart.

### C. Histogram of intervals forMIT-BIH Arrhythmia Database
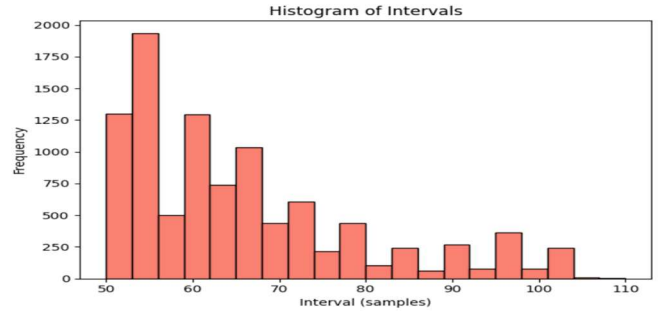

Fig 11: Histogram of intervals

Shows the distribution of R-R intervals (measured in samples).Main concentration between 50-70 samples, representing normal sinus rhythm, long tail extending to 110 samples. Multiple smaller peaks could indicate the premature beats (shorter intervals),compensatory pauses (longer intervals),different heart rate variations, possible rhythm abnormalities.

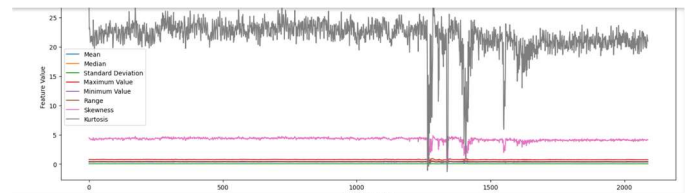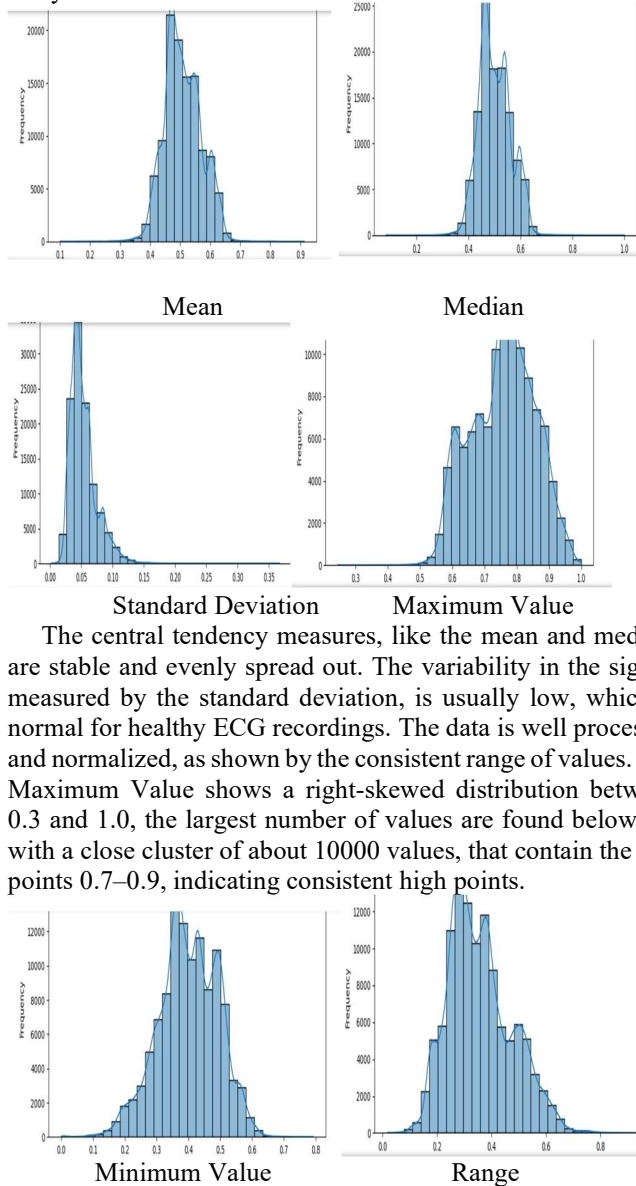### D. Feature plot for record 4


Fig 12: Feature Plot on Record 4

The Feature Plot on Record 4 displays the core descriptive statistics of the ECG signals. The kurtosis measurement, shown by the gray line, varies between 20 and 25, which means there is a lot of change in the recording. The skewness, represented by the pink line, stays steady around 4 to 5, indicating that the signal distribution is consistently uneven. Basic statistics like the mean (blue line), median (orange line), and standard deviation (green line) remain stable for most of the recording. However, there are noticeable disruptions around window 1500, pointing to changes in the signal that could be linked to arrhythmias.



Mean                                    Median



Standard Deviation          Maximum Value

The central tendency measures, like the mean and median, are stable and evenly spread out. The variability in the signal, measured by the standard deviation, is usually low, which is normal for healthy ECG recordings. The data is well processed and normalized, as shown by the consistent range of values. The Maximum Value shows a right-skewed distribution between 0.3 and 1.0, the largest number of values are found below 0.9 with a close cluster of about 10000 values, that contain the two points 0.7–0.9, indicating consistent high points.



Minimum Value                      Range
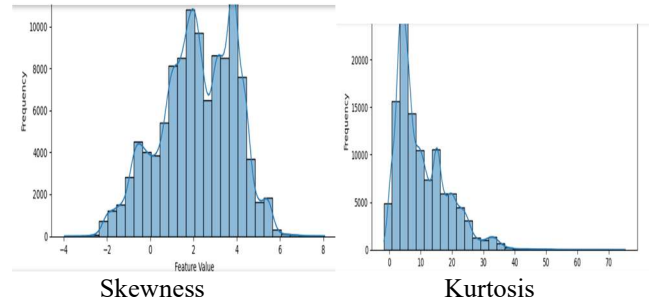


Skewness                                Kurtosis

Fig 13: Features of MIT-BIH Arrhythmia Database

The minimum value is between 0.4 and 0.5, equating to no rise or increase since very low frequencies of about 12000,
The Range is between 0.3 and 0.4, showing moderate variability from beat to beat. Skewness varies from -4 to 8, with peaks near 2 and 4, indicating different waveform shapes related to heart events. Kurtosis is mostly right-skewed with high kurtosis values lying above 5 reaching about twenty thousand, meaning, the distribution has mostly normal peaks but also feature some induced sharp or abnormal peaks.

## VII. DATA VISUALIZATION FOR P-WAVE

### A. Sliding Windows

Addresses a single ECG record (with the identification number ECG record 100). For this specific record numbered 100 it is found that 717 samples as the most appropriate window size which translates to about 1.99 based on the average distance between P-waves. This window size helps capture important patterns in each part of the ECG signal.
To carry out the analysis it is necessary to subdivide the signal into activities or actions, allowing for the creation of the 3,630 windows overall. The identification of P-waves is then undertaken in each segment. In the below figure there are left and right plots

- Left Plot: This contains the original ECG signal, raw and free from pre-processing with P-waves, marked with red dots. These red marks do the segmentation work. There is some small variation in amplitude in the raw data.

- Right Plot: This graph depicts the processed ECG signal, in which the noise is removed, and the essential information is available.

```
Analyzing record: 100
Optimal window size: 717 samples (1.99 seconds)
Number of segmented windows: 3630
```
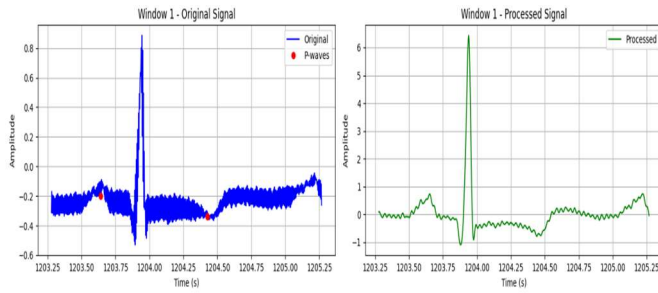
Fig :14 Sliding windows
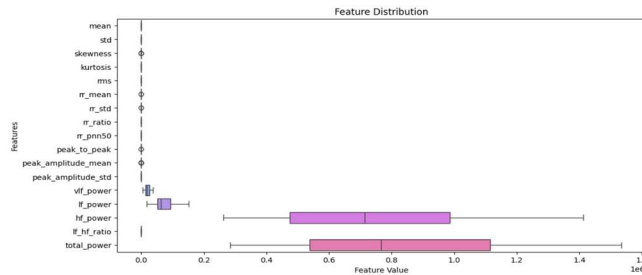
## B. Feature Distribution


Fig 15: Feature distribution

This box plot shows the spread of different signal analysis features. On the x-axis, values range from 0 to 1.6 million. One key takeaway is the large difference in scale between the features.Total Power, which includes the high frequency Power around 400,000 and 1.4 million, and has its medians at approximately 800,000. The low-frequency (lf) power has a more moderate range, centered near 100,000. Very low-frequency (vlf) power has a much narrower range close to zero.

Other features such as the mean, standard variations, and skewing factors, and peak amplitudes, etc., are represented much near zero almost, acting as fine-tuning parameters.As for RR interval features (rr_mean, rr_std, rr_ratio, rr_pnn50), show minimal variation, indicating their role in providing rhythm context and validating P-wave detection based on expected cardiac timing. This tight distribution supports accurate detection of heart rhythm.

## VIII. SPLITTING DATASET AND FEATURE SELECTION

### A. MIT-BIH Arrhythmia Database

Splitting dataset into 80% as training data and 20% as testing data. Before splitting data prepared data for classification, by performing label encoding to convert text labels into numeric values."

```
Number of classes: 15
Alphabetic classes: ['A' 'E' 'F' 'J' 'L' 'N' 'Q' 'R' 'S' 'V' 'a' 'e' 'f' 'j' 'x']
```

Using RFE (Recursive Feature Elimination) for feature selection. Selected the top 6 features.

```
Performing RFE to select top 6 features...
Selected features: [0 2 4 5 6 7]
```

### B. MIT-BIH Arrhythmia Database P-Wave Annotations

Splitting dataset into 80% as training data and 20% as testing data. Here also the RFE feature selection is used. There are 17 total features but by using recursive feature elimination selected the top 10 features for labelling with rr_std

```
Performing feature selection...

Selected features: ['mean' 'std' 'skewness' 'kurtosis' 'rms' 'rr_mean' 'rr_ratio' 'rr_pnn50'
 'peak_amplitude_std' 'hf_power']
```
Fig 16: Feature selection

## IX. MODELING

### A. MIT-BIH Arrhythmia Database

Multi-Class classification performed using Random Forest, Linear SVM, and Neural Network are trained and evaluated using randomized search cross-validation for hyperparameter optimization. Then, computes performance metrics such as accuracy, precision, and recall, displaying the best parameters for each model.

```
Random Forest Performance:
Best Parameters: {'n_estimators': 100, 'min_samples_split': 2, 'max_depth': 20}
```
Fig 17: Random forest best parameters

```
Linear SVM Performance:
Best Parameters: {'max_iter': 1000, 'dual': False, 'C': 1.0}
```
Fig 18: Linear SVM

```
Neural Network Performance:
Best Parameters: {'n_iter_no_change': 10, 'max_iter': 300, 'learning_rate_init': 0.001,
'hidden_layer_sizes': (100,), 'early_stopping': True, 'activation': 'relu'}
```
Fig 19: Neural network best parameters

```
Model Comparison:
            Model  Accuracy  Precision  Recall
0   Random Forest    0.9000     0.8909  0.9000
1      Linear SVM    0.7340     0.5997  0.7340
2  Neural Network    0.8839     0.8711  0.8839

Best performing model: Random Forest
Best model accuracy: 0.9000
```
Fig 20: Model comparison

The highest performance was exhibited by Random Forest model in comparison to the other two, with an accuracy level of 90%.Linear SVM model was rated the lowest as it is shortened in areas where the other models were doing better. The precision of the model was the most affected, which means it may not be very reliable for predicting positive cases. As for the Neural Network, it was not so perfect but achieved an approximate efficiency comparable to the random forest model.

### B. MIT-BIH Arrhythmia Database P-Wave Annotations

For binary classification, models are used for random forest, support vector machine , multi-layer perceptron neural network, and an unsupervised gaussian mixture model . Each supervised model is optimized using randomized search cross-validation. The unsupervised GMM model, in contrast, uses grid search cross-validation. rr_std, peak_to_peak, and hf_power are each used separately to create labels. If I select

rr_std as the y label, remove this feature from X, and each feature falls into a different category: rr_std is a time-domain feature, peak_to_peak is an amplitude feature, and hf_power is a frequency-domain feature. The performance metrics, like accuracy, precision, and recall, are calculated separately for each of these three features.

```
Model            Accuracy   Precision   Recall
RandomForest     1.000      1.000       1.000
SVM              1.000      1.000       1.000
Neural Network   0.667      1.000       0.500
GMM              0.667      1.000       0.500
```

Fig 21: performance metrics for rr_std labelling

```
Summary of Model Performance:

Model            Accuracy   Precision   Recall
RandomForest     0.333      0.500       0.500
SVM              0.333      0.500       0.500
Neural Network   0.333      0.500       0.500
GMM              0.667      1.000       0.500
```

Fig 22: performance metrics for peak_to_peak

```
Model            Accuracy   Precision   Recall
RandomForest     1.000      1.000       1.000
SVM              1.000      1.000       1.000
Neural Network   1.000      1.000       1.000
GMM              1.000      1.000       1.000
```

Fig 23: performance metrics for hf_power

## X. DISCUSSION

The Random Forest model showed even better results than the rest, with the difference, from 90.00%, precision 89.09% and Recall 90.00% degrees of precision. The linear support vector machine lagged, achieving 73.40% accuracy and 59.97% precision, indicating many false positives. The neural network possessed a total of 88.39% of recognition accuracy, which is quite good, which is however a smaller accuracy when compared to the random forests, giving more room for further research on enhancing its performance. Overall, the results indicate that the random forest model is the best in completing the classification assignment and the potential areas of improvement in terms of the neural network are foreseen. Random Forest and Support Vector Machine as their accuracy, precision and recall all achieved a perfect score. This could also suggest the capability of being able to classify P-wave annotations at high accuracy to perform the rr_std labeling. the neural network and gaussian mixture model , however, showed an overall accuracy of 66.7 percent. The issue of low recall rates in the Neural Network and GMM models should be addressed in future investigations by optimization of hyperparameters and by designing new features. When Y as hf_power then all four models got 100% of recall, accuracy and precision.

## XI. CONCLUSION

Random Forest model has a higher degree of accuracy, precision, and recall, the Neural Network has not peaked yet. The constraints of this research and the failure of required precision lack reveal the necessity for more advanced tools such as ensemble strategies including boosting and bagging, and kernel tricks of support vector machines. In so doing, there is a need to more adequately train the Neural Network and correct the errors for arrhythmia identification to be more accurate and less risk sensitive. The evaluation of P-wave annotations in the MIT-BIH Arrhythmia Database shows that Random Forest and Support Vector Machine techniques are currently the best for defining P-wave boundaries, given their high accuracy, precision, and recall. This indicates they are optimal for detecting arrhythmia conditions. In contrast, Neural Networks and GMMs still have issues which need to be addressed to make them sensitive and better perform. Waiting to be done are the application of further modifications while seeking on such ones as network parameter tuning and generation of new factors, which in this case are mainly resolution of their weaknesses, especially increments in recall.

## REFERENCE

[1] Qi M, Shao H, Shi N, Wang G, Lv Y. Arrhythmia classification detection based on multiple electrocardiograms databases. PLoS One. 2023 Sep 27;18(9):e0290995. doi: 10.1371/journal.pone.0290995. PMID: 37756278; PMCID: PMC10529562.

[2] Saclova, L., Nemcova, A., Smisek, R. *et al.* Reliable P wave detection in pathological ECG signals. *Sci Rep* **12**, 6589 (2022). https://doi.org/10.1038/s41598-022-10656-4