# Arrhythmia Detection Using ECG

## Important Features for Classification

Zach Reynolds
CS 390W
University of Massachusetts Amherst
Amherst, MA
zreynolds@umass.edu

John Rand
CS 390W
University of Massachusetts Amherst
Amherst, MA
johnrand@umass.edu

*Abstract*— **Early detection of cardiac arrhythmias is important in the health of any patient, emergency or otherwise, the following paper discusses the use of basic machine learning to determine important features of Electrocardiogram (ECG) data that are most influential when determining if a patient is at risk of having or should be diagnosed with a cardiac arrhythmia. Random Forest Classification is used to determine important features of ECG data with T-Tests used to determine difference in means between a healthy and unhealthy patient.**

*Keywords—arrhythmia, machine learning, random forest, t-test, ECG*

## I. INTRODUCTION

### A. Background

At its most basic description, cardiac arrhythmia is an irregular heartbeat. Cardiac arrhythmia occurs when there is a problem transmitting electrical signals that control the heart's beat. This can lead to a fast heart rate (tachycardia), a slow heart rate (bradycardia), or some form of irregular heartbeat. Within these types of arrhythmias, are several variations that are classified differently based on the pattern of the heartbeat, or the electrical activity in the heart. Arrhythmia is not always dangerous, but can be life threatening in situations if not appropriately treated. Symptoms can include chest pain, shortness of breath, sweating, syncope (fainting), and what is described as "a sense of impending doom."

Heartbeats are monitored through cardiac monitoring. Cardiac monitors, such as the Zoll X Series or Lifepak 15 cardiac monitors capture these heart trends through the use of a 12-lead Electrocardiogram. This is done by placing 10 conductive electrodes on the body to capture the electrical activity of the heart. Cardiac monitoring is done by trained healthcare professionals. This could be nurses in hospitals, or paramedics in a prehospital setting.

### B. Motivation

One of our researchers is an Emergency Medical Technician (EMT) who is familiar with the importance of cardiac monitoring in medical emergencies. EMTs in most states cannot interpret cardiac monitors on their own, however, some states allow EMTs to place a patient on cardiac monitoring and relay the information it outputs. As research in the field of cardiac monitoring continues, EMTs may be able

to better assess patients and have predictive algorithms built into the cardiac monitors to support them in their differential diagnosis.

## II. METHODOLOGY

The dataset was first loaded into Jupyter and analyzed for any missing data. It was found that the J Wave column had multiple indices of missing values. As a result, that column was dropped from use in our analysis.

### A. Random Forest Classification

Random Forest Classification was used in order to determine what features are important in determining if a patient has Arrhythmia. Random Forest Classification works by creating several decision trees that are then merged together to find an important value for each feature. These scores are from 0 to 1 and total to 1. The higher the value of the feature is, the more important it is in determining Arrhythmia. To do this, the data is split into training and testing data and is trained on the class of arrhythmia. For this analysis, we used a 70-30 training, testing split of the data. Which is standard in most cases.

After an initial classification was created, the feature importance values were viewed, and the 5 most important features were kept. These features were then trained again on another classifier to see if there was any improvement to accuracy. Metrics were taken on both classifiers to check precision, recall, and F1 score. A confusion matrix was then created and transferred onto a heatmap to visualize where the algorithm made Type I and Type II errors.

### B. Welch's T-Test

Welch's T-Test was used to verify the significance of the key attributes found from the random forest classification. Welch's T-Test compares the means of two populations with unequal variance. In this case, the two populations were people with an arrhythmia, and without. Five tests in total were run, with a different key attribute each time. The key attributes that were used for each test all reported back that the means of each population differed. This showed that each of the five attributes from the random forest classification proved to be significant indicators of someone with an arrhythmia. Along with each test being run, a box and whiskers plot of

each attribute with the control and population with arrhythmia was generated.
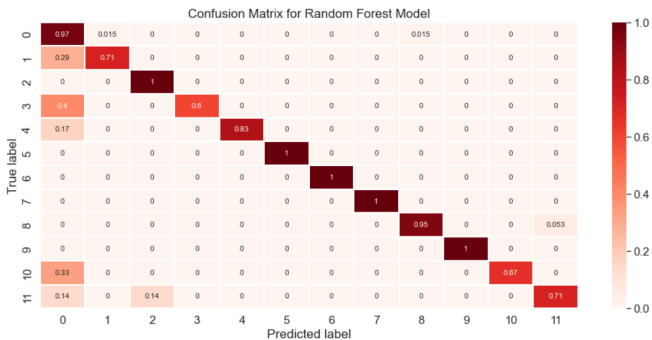
## III. RESULTS

The results gathered from this analysis were the results of the most accurate classifier, the metrics of the classifier, and t-tests performed on the 5 most important features determined by feature importance of the most accurate classifier.

### A. Random Forest

The results of the random forest classification revealed 5 main features that were most important in classifying Arrhythmia. They are as follows (in order of most important descending):

- Heart Rate
- V4 Node T-Wave Amplitude
- V5 Node QRSTA (Area of QRS segments + half the width of the T wave)
- V1 Node Number of Intrinsic Reflections
- AVR Node T-Wave Amplitude

Initial accuracy for the first classifier was at 70% with all features included. After adjusting the model to train on only these 5 features, accuracy improved to upwards of 90% in some iterations of the classifier. The following image is a heatmap created from the confusion matrix of our classifier predicting on the test data:



This heat map indicates that most of our misclassifications occurred when classifying class 0, which is considered a normal heart rhythm in the data. This is likely due to issues with the dataset which is further discussed in section IV. Below is the metrics discerned from the best fit classifier:

```
Accuracy: 0.9044117647058824
              precision    recall  f1-score   support

           1       0.88      0.97      0.92        66
           2       0.91      0.71      0.80        14
           3       0.83      1.00      0.91         5
           4       1.00      0.60      0.75         5
           5       1.00      0.83      0.91         6
           6       1.00      1.00      1.00         7
           8       1.00      1.00      1.00         1
           9       1.00      1.00      1.00         1
          10       0.95      0.95      0.95        19
          14       1.00      1.00      1.00         2
          15       1.00      0.67      0.80         3
          16       0.83      0.71      0.77         7

    accuracy                           0.90       136
   macro avg       0.95      0.87      0.90       136
weighted avg       0.91      0.90      0.90       136
```
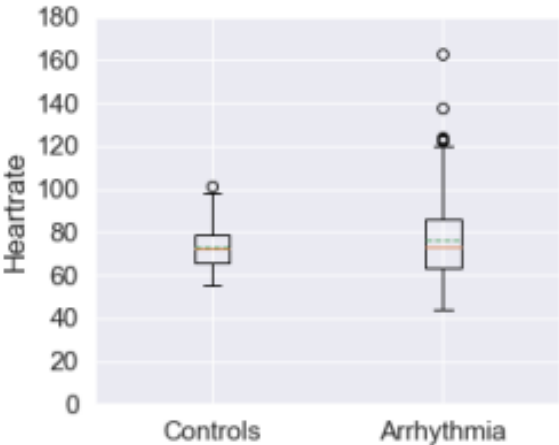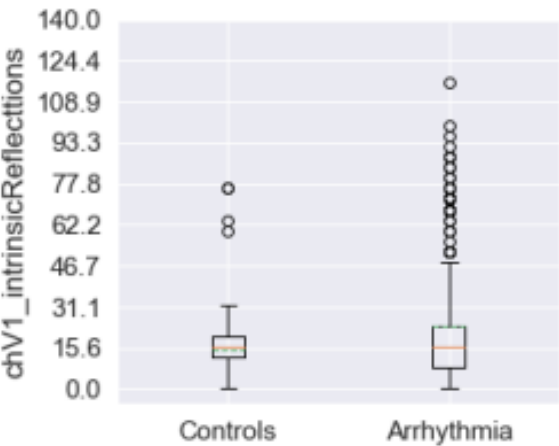
### B. Welch's T-Test

Null hypothesis for Welch's T-Test: The means of each attribute above were the same. Alternative hypothesis for Welch's T-Test: The means of each attribute above differed significantly.

Each test ran on the variables came back with a p-value < 0.05, in which we would reject the null hypothesis, stating that the means of each attribute between the two populations differed significantly.

Heart rate of control group and group with an arrhythmia.



V1 node number of intrinsic reflections for control group and group with an arrhythmia.



## IV. DISCUSSION

The results of the analysis make sense in regards to the circumstances of the analysis.

### A. Random Forest

Heart rate having the highest importance with regards to classification of Arrhythmia makes sense in the context because, as previously stated, Arrhythmia in its simplest form

is an abnormality in heart rate. The remaining high importance features were changes in electrical impulses, which also is correct in the context as electrical functions in the heart are meant to be consistent without significant changes.

The classifier also made misclassifications of Normal heart rates. This does, however, make sense in the context of the data because over half of the data (245 indices) were classified as class 0, which is the value assigned to normal heart rate. Because there are significantly more class 0 cases than any other cases, it makes sense that the data would be more likely to misclassify class 0, as it is the most common class in the testing set. This is also why some classes did not show up in the testing set due to low frequency.

*B. Welch's T-Test*

Running Welch's T-Test on the five attributes found by the random forest gave back a significant P-value for each attribute, in which the null hypothesis could be rejected, stating that the means of the control group and group with an arrhythmia differed. Showing that abnormal readings in these five classification attributes could be the onset of someone with an arrhythmia.

Taking a glance at the graphs, the control group has data that is mostly centered around the mean, with few outliers, whereas the group with an arrhythmia has many outliers and a different mean comparatively. This is the case for all five attribute graphs.

*C. Scores*

Values from Welch's T-Test:
Heart rate: T-value = -2.01 P-value = 4.50e-02
ChV4_TwaveAmp: T-value = 2.46 P-value = 1.43e-02
ChV5_QRSTA: T-value = 3.99 P-value = 8.03e-05

ChV1_intrinsicReflection: T-value: -5.08 P-value = 7.19e-07
ChAVR_TwaveAmp: T-value = -6.29 P-value = 9.98e-10

### Conclusion

Working with the UCI data allowed for some insights into important features that help when classifying ECGs from cardiac monitoring. However, the data is small in regards to this research. With so many classifications of arrhythmia, the data in the testing set often did not have a case of every class, which made it difficult to determine how well the classifier was working when many of the cases were class 0. This also explains why class 0 was most frequently misclassified because it had more than half of the cases in most testing sets that ran.

The T-Tests further validated our results as we accepted the alternative hypothesis of all T-Tests, which states there is a difference between average values in the controls and cases with arrhythmia. This likely implies that their values, whether higher or lower, have an effect on the heart's rhythm.

[1] *UCI Machine Learning Repository: Arrhythmia Data Set*. [Online]. Available: https://archive.ics.uci.edu/ml/datasets/Arrhythmia. [Accessed: 03-May-2022].

[2] "Heart arrhythmia," *Mayo Clinic*, 30-Apr-2022. [Online]. Available: https://www.mayoclinic.org/diseases-conditions/heart-arrhythmia/symptoms-causes/syc-20350668#:~:text=A%20heart%20arrhythmia%20(uh%2DRITH,slow%20(bradycardia)%20or%20irregularly. [Accessed: 03-May-2022].

[3] "Arrhythmia | irregular heartbeat," *MedlinePlus*. [Online]. Available: https://medlineplus.gov/arrhythmia.html. [Accessed: 04-May-2022].