

---

# Heart Disease Prediction from heart beat audio signals using Machine Learning

---

Ligari D. • Alberti A. <sup>1</sup>

<sup>1</sup> *Department of Computer Engineering, Data Science, University of Pavia, Italy*

*Course of Advanced Biomedical Machine Learning*

Github page: <https://github.com/DavideLigari01/advanced-biomedical-project>

Date: June 15, 2024

---

**Abstract** — —TO BE DEFINED—

**Keywords** — —TO BE DEFINED—

---

## CONTENTS

<b>1</b>	<b>Introduction</b>	<b>1</b>	<b>8</b>	<b>Models</b>	<b>9</b>
			8.1	Models Overview . . . . .	9
			8.1.1	LightGBM . . . . .	9
			8.1.2	XGBoost . . . . .	9
			8.1.3	CatBoost . . . . .	10
			8.1.4	Random Forest . . . . .	10
			8.1.5	MLP (Multilayer Perceptron) . . .	10
			8.2	Multiclass Model . . . . .	10
			8.2.1	Results . . . . .	10
			8.2.2	Risk Score . . . . .	10
			8.3	Normal VS Abnormal VS Artifact Model .	10
			8.3.1	Methodology . . . . .	10
			8.3.2	Results . . . . .	10
			8.3.3	Best Model Analysis . . . . .	11
			8.4	Diseases Only Model . . . . .	11
			8.4.1	Results . . . . .	11
<b>4</b>	<b>Goals</b>	<b>4</b>	<b>9</b>	<b>Analysis of the Best Model</b>	<b>11</b>
<b>5</b>	<b>Feature Extraction</b>	<b>4</b>	<b>10</b>	<b>Conclusion</b>	<b>11</b>
5.1	Features Type . . . . .	4	<b>11</b>	<b>Future Works</b>	<b>11</b>
5.2	Methodology . . . . .	5	<b>12</b>	<b>Appendix</b>	<b>12</b>
5.3	Sampling Rate and Interval Selection . . .	5			
5.3.1	Findings on Sampling Rate . . . . .	5			
5.3.2	Findings on Extraction Interval . . .	5			
5.4	Results . . . . .	6			
<b>6</b>	<b>Feature Selection</b>	<b>6</b>			
6.1	Determining the Optimal Number of Features	6			
6.2	Results and Analysis . . . . .	6			
6.3	Optimal Number of Features . . . . .	6			
<b>7</b>	<b>Correlation Analysis</b>	<b>6</b>			
7.1	Threshold Selection and Model Evaluation .	8			
7.2	Selected Features and Correlation Matrix .	8			

## 1. INTRODUCTION

Heart disease remains a leading cause of mortality world-wide. Early diagnosis is critical for effective treatment and management. Traditional methods of diagnosis often involve invasive procedures and expensive equipment. Recent advancements in machine learning have opened new avenues for non-invasive diagnosis using heart sound recordings. This paper explores the use of machine learning and network analysis to predict heart disease from heart beat audio signals. Despite significant progress, gaps remain in accurately classifying heart sounds due to data

imbalance and the presence of noise in recordings. This study aims to address these gaps by employing advanced data preprocessing techniques and robust machine learning models. The research question guiding this study is: How can machine learning models be optimized to improve the accuracy of heart disease prediction from heart sound recordings?

## 2. DATASET DESCRIPTION

The dataset for this project was obtained from a Kaggle repository titled *Dangerous Heartbeat Dataset (DHD)* [1], which in turn sources its data from the PASCAL Classifying Heart Sounds Challenge 2011 (CHSC2011) [2]. This dataset comprises audio recordings of heartbeats, categorized into different types of heart sounds. Specifically, the dataset consists of 5 types of recordings: Normal Heart Sounds, Murmur Sounds, Extra Heart Sounds, Extrasystole Sounds, and Artifacts. Data has been gathered from the general public via the iStethoscope Pro iPhone app and from a clinic trial in hospitals using the digital stethoscope DigiScope.

### 2.1. Data Distribution

Figure 1, show the presence of highly imbalanced classes in the dataset. This poses a challenge for the classification task as the model may not have enough samples to learn from, especially for the 'Extrastole' and 'Extrahls' classes. This problem is tackled trying to augment the data available, both by segmenting the audio files and by using data augmentation techniques. Furthermore, we test the effectiveness of oversampling and undersampling techniques on the model performance. The data is split into training and testing sets, with a 80% - 20% ratio, respectively. The validation set is omitted, due to the low number of samples available.

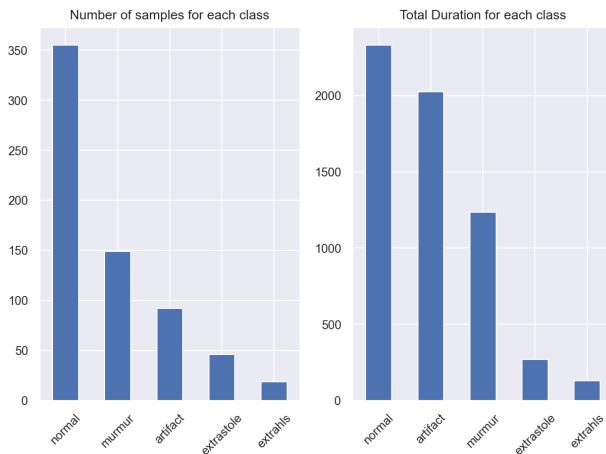


Fig. 1: Number of samples per duration.

### 2.2. Data Preprocessing

Here is a list of the preprocessing operations performed on the dataset:

**Noise Reduction:** the audio data was already provided in

a clipped format to minimize noise and irrelevant information.

**Removal of Corrupted Files:** corrupted files were identified and removed from the dataset to ensure data quality.

**Resampling:** all audio files were resampled to a common frequency of 4000 Hz, which was identified as the optimal sampling rate (see Section 5.3).

**Segmentation:** the audio data was segmented into 1-second intervals, identified as the optimal extraction interval (see Section 5.3).

**Outlier Detection and Removal:** we investigated the average duration of each class and found that the 'artifact' class had a significantly larger average duration. This was due to a few lengthy audio recordings (see Figure 2). These recordings were considered as outliers and removed from the dataset, as a large number of samples from the same audio might not be as informative.

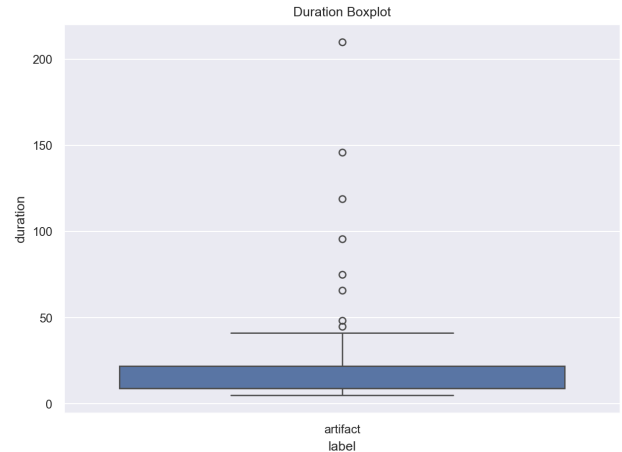


Fig. 2: Outliers in the Artifacts class.

## 3. HEART SOUND CATEGORIES

Heart sounds can be categorized into different classes based on their acoustic characteristics and clinical significance. Accurate classification of these sounds is essential for diagnosing and treating a variety of cardiac conditions. The primary categories include Normal heart sounds, Murmurs, Extra Heart Sounds, Artifacts, and Extra Systoles. Understanding the distinct features and clinical implications of each class is a crucial step before building a machine learning model to classify heartbeats. This phase is particularly important for the identification of patterns that are characteristic of specific classes, which in turn guides the selection of features to extract from the audio. This knowledge aids in identifying specific patterns and anomalies within the heart sounds, leading to more precise and reliable model predictions.

### 3.1. Normal

The Normal category includes recordings of typical, healthy heart sounds. These sounds exhibit the characteris-

tic “lub-dub, lub-dub” pattern, where “lub” (S1) represents the closing of the atrioventricular valves and “dub” (S2) signifies the closing of the semilunar valves. In a normal heart, the time interval between “lub” and “dub” is shorter than the interval from “dub” to the next “lub,” especially when the heart rate is below 140 beats per minute. Most normal heart rates at rest fall between 60 and 100 beats per minute, though rates can vary from 40 to 140 beats per minute based on factors such as age and activity level. Recordings may include background noises like traffic or radio sounds and may capture incidental noises such as breathing or microphone contact with clothing or skin. It contains both clean and noisy normal recordings, the latter featuring significant background noise or distortion, which simulates real-world conditions.

Figure 3 shows a sample of a normal heart beat audio. The characteristic “lub-dub, lub-dub” pattern can be observed, where the peaks represent the “lub” (S1) and “dub” (S2) sounds of a healthy heart.

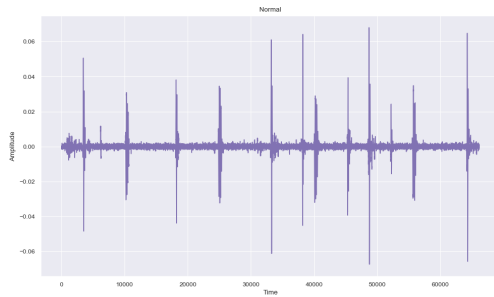


Fig. 3: Sample of normal heart beat audio.

### 3.2. Murmur

Heart murmurs are abnormal sounds during the heartbeat cycle, such as a “whooshing, roaring, rumbling, or turbulent fluid” noise, heard between the “lub” and “dub” (systolic murmur) or between “dub” and “lub” (diastolic murmur). These murmurs are typically indicative of turbulent blood flow in the heart and can signal various heart conditions, some of which may be serious. It is crucial to distinguish murmurs from the normal “lub-dub” sounds since they occur between the primary heart sounds and not concurrently with them. It also includes noisy murmur data, which mimics real-world recording scenarios by incorporating significant background noise and distortions.

Figure 4 shows a sample of a murmur heart beat audio. The presence of additional sounds between the “lub” and “dub” peaks can be observed, indicating the characteristic “whooshing, roaring, rumbling, or turbulent fluid” noise typical of heart murmurs.

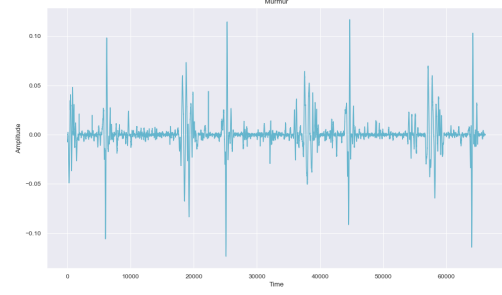


Fig. 4: Sample of murmur heart beat audio.

### 3.3. Extra Heart Sound

Extra heart sounds are characterized by an additional sound in the cardiac cycle, producing patterns such as “lub-lub dub” or “lub dub-dub”. These sounds can arise from physiological or pathological conditions. For example, a third heart sound (S3) may indicate heart failure or volume overload, while a fourth heart sound (S4) can be associated with a stiff or hypertrophic ventricle. Detecting these extra sounds is important for identifying potential heart diseases early, allowing for timely intervention and management. Figure 5 shows a sample of an extra heart sound audio. The presence of additional peaks within the normal “lub-dub” pattern indicates extra heart sounds, which can be critical for diagnosing various heart conditions.

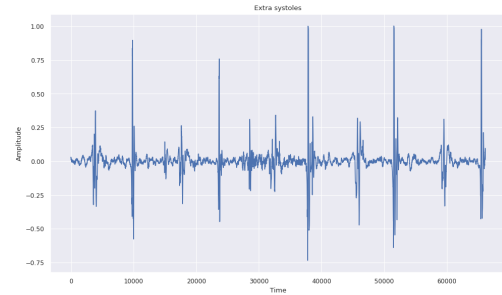


Fig. 5: Sample of extra heart sound audio.

### 3.4. Artifact

The Artifact category consists of recordings with non-cardiac sounds, including feedback squeals, echoes, speech, music, and various types of noise. These recordings generally lack discernible heart sounds and do not exhibit the temporal periodicity typical of heartbeats at frequencies below 195 Hz. Accurately identifying artifacts is essential to avoid misinterpreting non-cardiac sounds as pathological heart sounds, ensuring that data collection efforts focus on genuine heart sounds. Figure 6 shows a sample of an artifact heart beat audio, there can be observed that there is not a clear pattern in the audio.

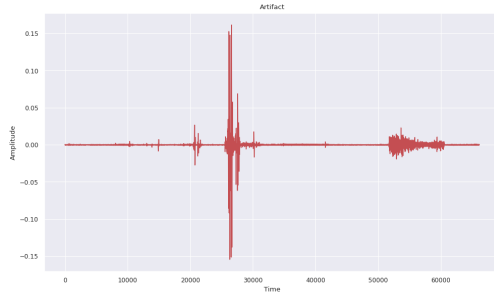


Fig. 6: Sample of artifact heart beat audio

### 3.5. Extra systoles

Extra systoles refers to extra or skipped heartbeats, resulting in irregular patterns such as “lub-lub dub” or “lub dub-dub”. Unlike the regular extra heart sounds, extra systoles are sporadic and do not follow a consistent rhythm. These premature beats can occur in healthy individuals, particularly children, but they may also be associated with various heart diseases. Identifying extra systoles is crucial as they can be early indicators of cardiac conditions that might require medical attention if they occur frequently or in certain patterns.

In the audio signal depicted in Figure 7, irregularities within the normal “lub-dub” pattern are evident. These irregularities manifest as additional peaks or skipped beats, indicating extra systoles.

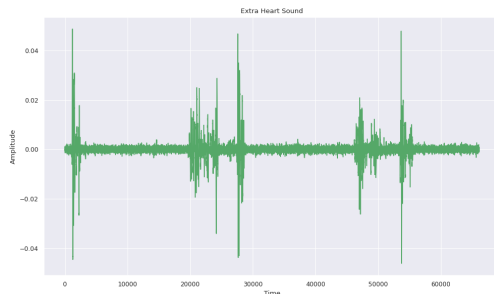


Fig. 7: Sample of extra systoles heart beat audio

### 3.6. Comparison of Heart Sounds

In Figure 8, a comparison of the different classes of heart sounds can be observed.

As we can see, the “Artifact” signal appears erratic with no consistent pattern, likely representing noise or interference rather than true heart sounds.

The “Murmurs” signal shows irregular fluctuations in amplitude, which could indicate turbulent blood flow typically associated with murmurs.

The signal for “Extra Heart Beat Sound” has occasional spikes in amplitude that stand out from the baseline.

The “Normal” signal appears more uniform and regular compared to the others, reflecting the expected rhythm of a healthy heartbeat. Finally, the signal for “Extra Systoles” shows extra spikes at irregular intervals, indicating unexpected contractions of the heart muscle (systoles) occurring outside the normal rhythm.

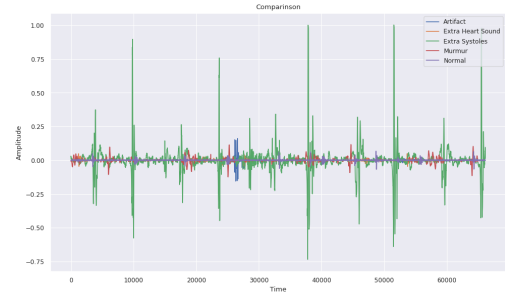


Fig. 8: Comparison of the different classes of heart sounds.

## 4. GOALS

## 5. FEATURE EXTRACTION

As demonstrated by [4] and [3], MFCCs are highly effective features for heartbeat classification. In addition to MFCCs, we incorporated other features to capture various characteristics of heart sounds, enhancing the classification accuracy. The features used are explained in the following section.

### 5.1. Features Type

#### MFCC

Mel-Frequency Cepstral Coefficients (MFCCs) are representations of the short-term power spectrum of sound. They are derived by taking the Fourier transform of a signal, mapping the powers of the spectrum onto the mel scale, taking the logarithm, and then performing a discrete cosine transform. MFCCs are effective in capturing the timbral texture of audio and are widely used in speech and audio processing due to their ability to represent the envelope of the time power spectrum. In heartbeat classification, MFCCs can reflect the different perceived quality of heart sounds, such as the presence of murmurs or other anomalies.

#### Chroma STFT

Chroma features represent the 12 different pitch classes of music (e.g., C, C#, D, etc.). They are particularly useful for capturing harmonic and melodic characteristics in music. By mapping audio signals onto the chroma scale, these features can identify pitches regardless of the octave, making them useful for analyzing harmonic content in heart sounds.

#### RMS

Root Mean Square (RMS) measures the magnitude of varying quantities, in this case, the amplitude of an audio signal. It is a straightforward way to compute the energy of the signal over a given time frame. RMS is useful in audio analysis for detecting volume changes and can help identify different types of heartbeats based on their energy levels. For example, in a given timeframe the RMS may be altered by the presence of a murmur with respect to a normal heart sound.

### ZCR

Zero-Crossing Rate (ZCR) is the rate at which a signal changes sign, indicating how often the signal crosses the zero amplitude line. It is particularly useful for detecting the noisiness and the temporal structure of the signal. In heartbeat classification, ZCR can help differentiate between normal and abnormal sounds by highlighting changes in signal periodicity.

### CQT

Constant-Q Transform (CQT) is a time-frequency representation with a logarithmic frequency scale, making it suitable for musical applications. Since it captures more detail at lower frequencies, it may be useful for analyzing the low-frequency components of heart sounds.

### Spectral Centroid

The spectral centroid indicates the center of mass of the spectrum and is often perceived as the brightness of a sound. It is calculated as the weighted mean of the frequencies present in the signal, with their magnitudes as weights. In heart sound analysis, a higher spectral centroid can indicate sharper, more pronounced sounds, while a lower centroid suggests smoother sounds.

### Spectral Bandwidth

Spectral bandwidth measures the width of the spectrum around the centroid, providing an indication of the range of frequencies present. It is calculated as the square root of the variance of the spectrum. This feature helps in understanding the spread of the frequency components in the heart sounds, which can be indicative of different heart conditions.

**Spectral Roll-off** Spectral roll-off is the frequency below which a certain percentage of the total spectral energy lies. It is typically set at 85% and helps distinguish between harmonic and non-harmonic content. In heartbeat classification, spectral roll-off can be used to differentiate between sounds with a concentrated energy distribution and those with more dispersed energy.

## 5.2. Methodology

The features were extracted from the audio signals using the *Librosa* library in Python. It is worth to underline four main aspects in the extraction methodology having a direct impact on the results:

- **Normalization:** the audio are loaded using the *torchaudio.load()* function, which normalized the audio signals in the range  $[-1, 1]$ . This is important to ensure that the features are on the same scale and to prevent the model from being biased towards features with larger values.
- **Audio Clipping:** to extract features, we divided the audio recordings into chunks of a given length (in seconds). This segmentation allowed us to analyze the impact of different extraction intervals on model performance, additionally it allow for augmenting the data available.

- **Sampling Rate:** while literature on spoken language often suggests that 16000 Hz is sufficient, it was necessary to assess the best sampling rate for heartbeat sounds specifically. We evaluated two sampling rates to determine the optimal rate for heartbeat sounds.

- **Hop and Window Size:** the hop size determines the number of samples between successive windows, while the window size determines the number of samples considered. Each feature was extracted using the same window length and hop length facilitating a fair assessment of each feature's contribution to the classification task.

The library used for the extraction is *Librosa*, which is a Python package for music and audio analysis.

## 5.3. Sampling Rate and Interval Selection

To determine the optimal sampling rate and extraction interval for heartbeat classification, a series of experiments was conducted. We trained various models (Random Forest, Support Vector Machine, Logistic Regression) with different features (Chroma, MFCC\_30, MFCC\_120, CQT\_30, CQT\_70), sampling rates (mixed, 4000 Hz), and extraction intervals (0.5s, 1s, 2s, 3s). Features were used as extracted without any additional processing. The models were trained on 80% of the data and tested on the remaining 20%.

### 5.3.1. Findings on Sampling Rate

Our experiments revealed no clear advantage of using a mix of frequencies over a fixed sample rate, independently of the metric used. Specifically, a fixed sampling rate of 4000 Hz was found to reduce the risk of bias introduction, improve efficiency, and enable the use of a wider variety of features. This sampling rate provided a consistent and reliable basis for feature extraction, as demonstrated in Figure 9.

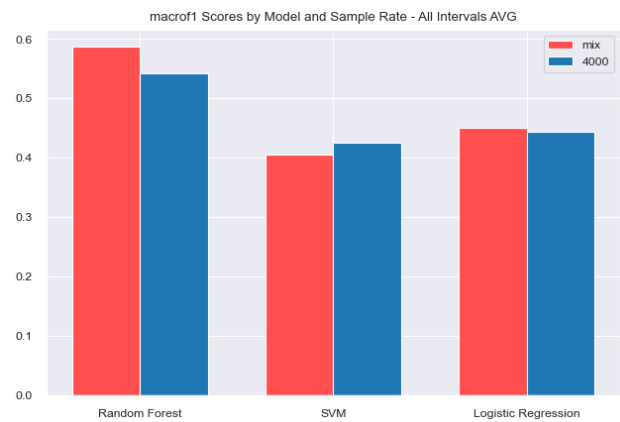


Fig. 9: Comparison of the macro F1 score for different sampling rates.

### 5.3.2. Findings on Extraction Interval

The choice of extraction interval had a significant impact on the number of samples and the class distribution, as shown in Figure 10. To address the high class imbalance,



we used the macro F1 score and Matthews Correlation Coefficient (MCC) as evaluation metrics. Our results indicated that a 2-second interval yielded the best performance. However, this choice also reduced the number of samples, potentially causing issues during training and testing, especially with more complex models.

As a compromise, we recommend using a 1-second interval. This interval offers a good balance between the number of samples and the class distribution, ensuring robust model performance while maintaining a sufficient dataset size. Figures 11 and 12 illustrate the impact of different extraction intervals on the macro F1 score and MCC, respectively.

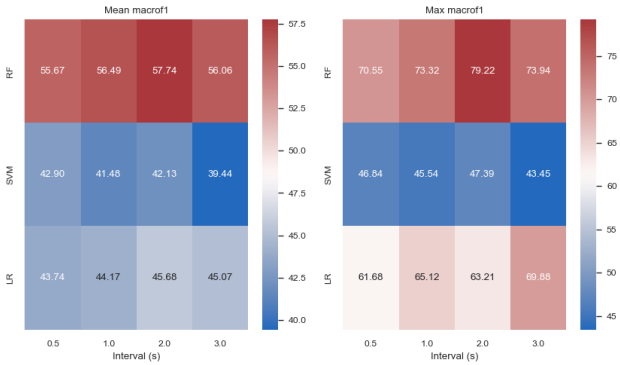


Fig. 11: Comparison of the macro F1 score for different extraction intervals.

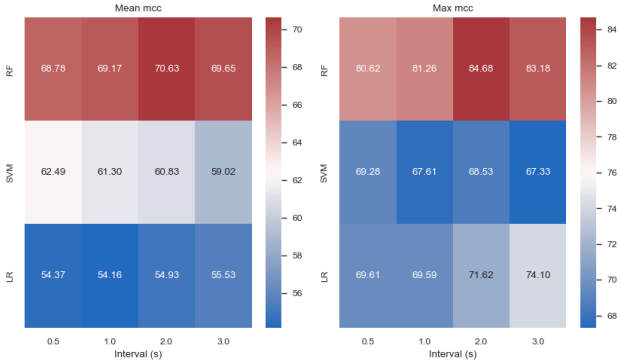


Fig. 12: Comparison of the MCC for different extraction intervals.

## 5.4. Results

## 6. FEATURE SELECTION

In this section, the focus is on identifying the optimal number of features for different types of audio features (MFCC, Chroma, CQT, etc.). Each type of feature can consist of a varying number of individual features, and it is crucial to determine the optimal number to maximize the model's performance. Proper feature selection is vital because it directly impacts the efficiency, accuracy, and generalizability of the machine learning model. Including too many features can lead to overfitting, increased computational costs, and degraded performance due to the curse of dimensionality. On the other hand, using too few features can result in underfitting, where the model fails to capture the necessary patterns in the data. Therefore, the goal of this study

is to identify the optimal number of features for each type to ensure the model is both robust and efficient.

### 6.1. Determining the Optimal Number of Features

To maximize the model's performance, it is crucial to determine the optimal number of features for each type. This was achieved using a 'One Model per Feature' approach, which involved the following steps:

#### Algorithm 1 Feature Optimization Process

- 1: **Step 1:** Extract feature sets with varying sizes: 12, 20, 30, 40, 60, 70, 90, and 120 features.
- 2: **Step 2:** Train a separate model for each feature type and feature set size.
- 3: **Step 3:** For each feature set, train three different models using three classifiers: SVM, Random Forest, and Logistic Regression.
- 4: **Step 4:** Evaluate the performance of each model.

Figure 13 shows the results obtained for each type of feature with the Random Forest model, which significantly outperformed the other classifiers.

### 6.2. Results and Analysis

**MFCC:** The MFCC features consistently achieved the highest F1 scores, peaking around 0.7. This indicates their effectiveness for the classification task. Interestingly, the number of MFCC features (ranging from 30 to 120) did not drastically affect performance, suggesting that even a smaller set of MFCC features can be highly informative.

**CQT:** The CQT features showed moderate performance, with F1 scores around 0.4 to 0.5. The optimal number of features was around 70, beyond which there was no significant improvement.

**RMS:** RMS features exhibited F1 scores ranging from 0.4 to 0.5, with optimal performance achieved with around 70 features.

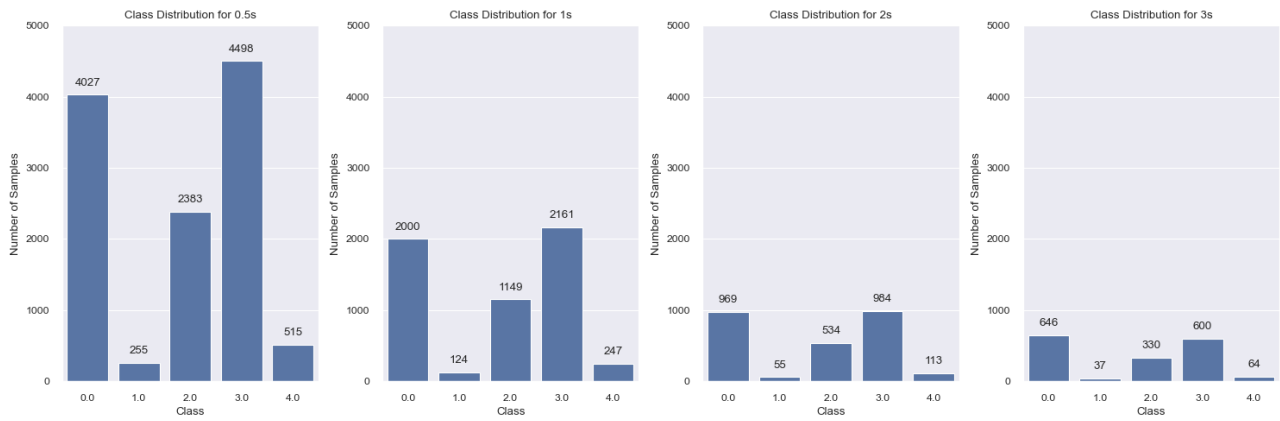
**ZCR, Spectral Centroid (SC), Spectral Bandwidth (SB), and Spectral Roll-off (SR):** The F1 scores for these features generally stabilized around 0.4 to 0.5. Increasing the number of features beyond 40 did not result in significant performance gains and could even degrade the model's performance. This suggests that adding too many features, especially those without strong predictive power, can confuse the model and degrade its performance.

### 6.3. Optimal Number of Features

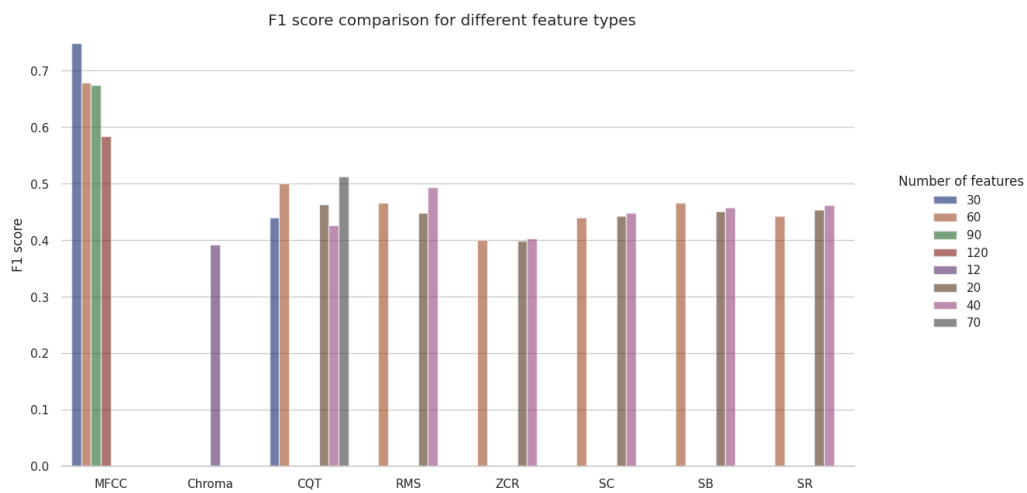
Based on the results, the optimal number of features for each type is shown in Table 1.

## 7. CORRELATION ANALYSIS

Given the large number of features (338 in total), it was necessary to identify and remove features that are poorly correlated with the target variable as well as those that are highly correlated with each other. Due to the high number of features, a visual approach, such as a correlation matrix,



**Fig. 10:** Class distribution for different extraction intervals.



**Fig. 13:** F1 score per number of features

Feature Type	Optimal Number of Features
MFCC	30
Chroma	12
CQT	70
RMS	40
ZCR	40
Spectral Centroid	40
Spectral Bandwidth	60
Spectral Roll-off	40

Table 1: Optimal number of features for each type

was not feasible. Instead, two filters were applied to select the most relevant features using the Spearman correlation coefficient, as the normality test failed.

The first filter is based on the correlation between the features and the target variable. Features with a correlation below a certain threshold with the target variable are removed.

The second filter focuses on the correlation among the features themselves. It counts, for each feature, the number of other features with which it has a correlation above a certain threshold. Features with a number of correlations above a specified threshold are then removed.

#### Algorithm 2 Feature Selection Process

- 1: **Step 1:** Compute the normal test (D'Agostino Pearson).
- 2: **Step 2:** Compute the Spearman correlation coefficient for each feature with the target variable.
- 3: **Step 3:** Apply the first filter to remove features with a correlation below a certain threshold with the target variable.
- 4: **Step 4:** Compute the correlation matrix among all features.
- 5: **Step 5:** Apply the second filter to remove features that have a high number of correlations (above a certain threshold) with other features.
- 6: **Step 6:** Choose threshold values empirically and apply the filters using various combinations of these thresholds.
- 7: **Step 7:** Train Random Forest models on the filtered data to evaluate performance and select the best combination of thresholds.

### 7.1. Threshold Selection and Model Evaluation

Threshold values were chosen empirically and the filters were applied using the combinations shown in Table 2. Using the filtered data, Random Forest models were trained and evaluated, as Random Forest was found to be the best performing model. The optimal combination of thresholds was found to be: threshold 1 = 0, threshold 2 = 0.6 and number of features = 30, resulting in 30 features.

With threshold 1 = 0, the filter on the correlation between the features and the target variable was effectively bypassed. However, with threshold 2 = 0.6, a stringent filter was applied on the correlation among the features them-

Threshold	Values
THRESHOLD 1	0 - 0.1 - 0.2 - 0.3 - 0.4 - 0.5
THRESHOLD 2	0.6 - 0.7 - 0.8 - 0.9 - 1
N° FEATURES	5 - 10 - 15 - 20 - 25 - 30 - 40

Table 2: threshold values

selves, removing features that had a correlation above 0.6 with at least 30 other features. This indicates that having features highly correlated with each other is more detrimental to the model than having features poorly correlated with the target variable.

Figure 14 shows the results obtained with the model trained on filtered features compared to the model trained on all features. As demonstrated, the model trained on filtered features performs significantly better.

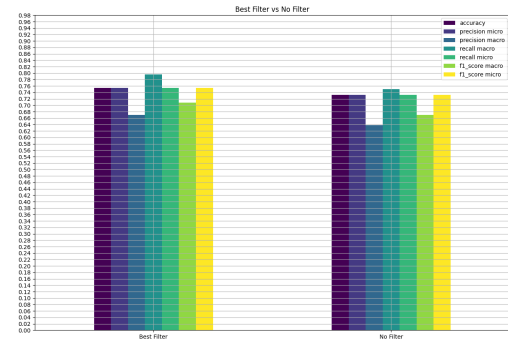


Fig. 14: Comparison of different metrics between the model on all features and the model on the filtered ones

### 7.2. Selected Features and Correlation Matrix

From this analysis, 41 features remained: 28 MFCC, 12 Chroma, and 1 ZCR. The correlation matrix of the filtered features is shown in Figure 15. This matrix illustrates the pairwise correlation between the selected features, with the color intensity indicating the strength and direction of the correlation. Dark red cells represent high positive correlations, while dark blue cells indicate high negative correlations.

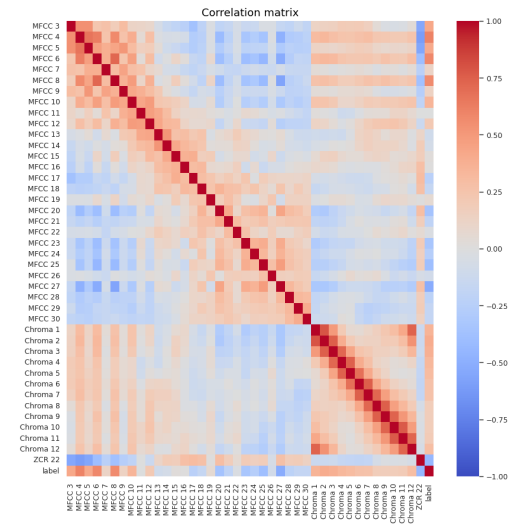


Fig. 15: Correlation matrix of the filtered features

The matrix demonstrates that the remaining features have



low correlations with each other, as evidenced by the predominantly light colors away from the diagonal. This implies that the features are relatively uncorrelated, preventing multicollinearity issues and enhancing the robustness of the model. The high diagonal values indicate that each feature is perfectly correlated with itself, which is expected. However, the off-diagonal values being close to zero for most feature pairs confirm that the filtering process was effective in selecting features that do not exhibit high inter-correlations.

## 8. MODELS

To effectively address the classification problem, we strategically split it into three distinct tasks:

- 1. Multiclass Classification:** This task addresses the broadest scope of prediction by classifying heartbeats into one of the five classes present in the dataset. This is the most challenging task due to the imbalanced nature of the classes and the need to differentiate between multiple disease categories.
- 2. Normal vs Abnormal vs Artifact Classification:** This task takes a more focused approach by categorizing heartbeats into three distinct classes: normal, abnormal, or artifact. It mitigates the challenge of imbalanced classes present in the dataset and allows for a better understanding of the nuances within each category. Additionally, it allows for a false positive rate minimization scenario, where the final goal is not the overall accuracy, but the minimization of false normal predictions. This is particularly important in a medical context, where misclassifying an abnormal heartbeat as normal can have severe consequences.
- 3. Diseases-only Classification:** This task narrows the focus exclusively to the prediction of diseases, excluding normal and artifact classes. It allows for a deeper exploration of the underlying diseases present in the dataset while filtering out irrelevant data points. By isolating disease-related patterns and features, the model can prioritize accuracy in disease prediction without the potential confounding factors introduced by normal and artifact heartbeats.

In all cases we trained a variety of models, with the aim of finding the best-performing model for each task. The models used in this study are listed in Table 3.

The working principle of the basics models is briefly explained in the next section. As regard the ensemble models, they are built by combining the predictions of multiple base models. Specifically, from MLP\_Ensemble1 to MLP\_Ensemble6, we combined the predictions of the models listed in the name. The combination was done using the *VotingClassifier* from the *sklearn* library, with the *soft* vote strategy. The ALL\_Ensemble model combines the predictions of all the models above it, using a majority vote strategy. The CB\_ALL\_Ensemble model is similar to the ALL\_Ensemble model, but it uses the CatBoost algorithm to combine the predictions of the base models.

Name	Architecture (hidden layers)
Random Forest	-
XGBoost	-
CatBoost	-
LightGBM	-
MLP_Basic	(128, 64, 32)
MLP_Ultra	(512, 256, 128, 64, 32)
MLP_Large	(256, 128, 64, 32)
MLP_Small	(64, 32)
MLP_Tiny	(32, 16)
MLP_Reverse	(32, 64, 128, 256, 512, 256, 128, 64, 32)
MLP_Bottleneck	(512, 64, 32)
MLP_Rollercoaster	(512, 128, 256, 128, 256, 64, 32)
MLP_Hourglass	(512, 256, 128, 64, 32, 64, 128, 256, 512)
MLP_Pyramid	(1024, 512, 256, 128, 128, 128, 64, 32)
MLP_Wide	(1024, 1024)
MLP_WideUltra	(1024, 1024, 128, 32)
MLP_Sparse	(32, 16, 8)
MLP_Dropout	(128, 64, 32)
MLP_Ensemble1	MLP_Basic, Large, Ultra
MLP_Ensemble2	RandomForest, MLP_Ultra
MLP_Ensemble3	MLP_Rollercoaster, Large
MLP_Ensemble4	MLP_Rollercoaster, Large, Ultra
MLP_Ensemble5	RandomForest, MLP_Ultra, Rollercoaster
MLP_Ensemble6	MLP_Rollercoaster, Large, Ultra, Wide
ALL_Ensemble	All models majority vote
CB_ALL_Ensemble	All models CatBoost

Table 3: Models names and architectures.

### 8.1. Models Overview

#### 8.1.1. LightGBM

LightGBM (Light Gradient Boosting Machine) is a gradient boosting framework that builds decision trees in a leaf-wise manner. Unlike level-wise growth used in other boosting algorithms, LightGBM grows trees leaf-wise, leading to better optimization.

The objective function to minimize is:

$$L = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_k \Omega(f_k)$$

where  $l$  is the loss function (e.g., mean squared error for regression),  $y_i$  is the true label,  $\hat{y}_i$  is the predicted label, and  $\Omega$  is the regularization term to avoid overfitting.

#### 8.1.2. XGBoost

XGBoost (Extreme Gradient Boosting) enhances the gradient boosting technique by optimizing the loss function using second-order Taylor expansion. The objective function is:

$$L(\theta) = \sum_{i=1}^n l(y_i, \hat{y}_i(t)) + \sum_k \Omega(f_k)$$

where  $l$  is the loss function,  $y_i$  is the true value,  $\hat{y}_i(t)$  is the prediction at the  $t$ -th iteration, and  $\Omega(f_k)$  is the regularization term. The regularization term is given by:

$$\Omega(f) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

where  $T$  is the number of leaves,  $w_j$  are the leaf weights, and  $\gamma$  and  $\lambda$  are regularization parameters.

### 8.1.3. CatBoost

CatBoost (Categorical Boosting) is designed to handle categorical features natively. It uses ordered boosting to avoid overfitting and provides efficient handling of categorical data without extensive preprocessing. The model learns using the following objective function:

$$L = \sum_{i=1}^n l(y_i, \hat{y}_i) + \lambda \sum_{j=1}^m ||w_j||^2$$

where  $l$  is the loss function,  $y_i$  is the true label,  $\hat{y}_i$  is the predicted label,  $\lambda$  is the regularization parameter, and  $w_j$  are the model weights.

### 8.1.4. Random Forest

Random Forest is an ensemble method that builds multiple decision trees using bootstrap samples and random feature subsets. The final prediction is made by averaging the predictions (regression) or majority voting (classification). The prediction for a regression problem is:

$$\hat{y} = \frac{1}{T} \sum_{t=1}^T h_t(x)$$

where  $T$  is the number of trees,  $h_t$  is the prediction of the  $t$ -th tree, and  $x$  is the input feature vector.

### 8.1.5. MLP (Multilayer Perceptron)

MLP (Multilayer Perceptron) is a type of feedforward neural network with one or more hidden layers. Each neuron computes a weighted sum of inputs, applies an activation function, and passes the result to the next layer. The output of a neuron in layer  $l$  is:

$$a_i^{(l)} = f \left( \sum_{j=1}^{n^{(l-1)}} w_{ij}^{(l)} a_j^{(l-1)} + b_i^{(l)} \right)$$

where  $a_i^{(l)}$  is the activation of the  $i$ -th neuron in layer  $l$ ,  $w_{ij}^{(l)}$  is the weight between neuron  $j$  in layer  $l-1$  and neuron  $i$  in layer  $l$ ,  $b_i^{(l)}$  is the bias term, and  $f$  is the activation function (e.g., ReLU, sigmoid).

## 8.2. Multiclass Model

### 8.2.1. Results

### 8.2.2. Risk Score

## 8.3. Normal VS Abnormal VS Artifact Model

One of the primary concerns in applying machine learning models to the medical field is the potential consequence of misclassifying abnormal conditions as normal. Specifically, in the case of heartbeats, misclassifying an abnormal heartbeat as normal can have severe clinical implications. Therefore, our focus shifts from achieving the highest overall accuracy to minimizing the false positive rate (FPR) for the normal class. This approach ensures that fewer abnormal heartbeats are incorrectly classified as normal, thereby enhancing patient safety and diagnostic reliability.

### 8.3.1. Methodology

To effectively address the classification challenge, we categorized the heartbeats into three classes: normal, abnormal, and artifact, with the latter encompassing extrasystoles, murmurs, and extrahls. We then selected models from the previous section, prioritizing those with a low risk score. The features utilized were the filtered ones, and no class balancing was necessary since the classes were inherently balanced.

Each model was trained on 80% of the dataset and tested on the remaining 20%. The evaluation metrics included the ROC curve, macro F1 score, and balanced accuracy.

### 8.3.2. Results

The initial evaluation is presented using the ROC curves for all models. To address the three-class problem, we employed the one-vs-rest strategy, comparing each class against the others. Figure 16 displays the ROC curves for the normal class versus the other classes across all models.

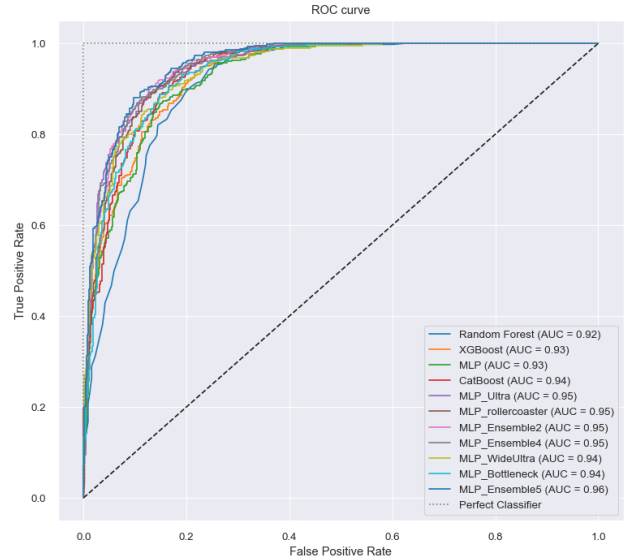


Fig. 16: ROC curves for the normal class against the rest of the classes across all models.

The results, indicated by the AUC values, show that the MLP model outperformed the others. Specifically, *MLP\_Ensemble5* achieved the highest AUC value of 0.96, followed by *MLP\_Ultra*, *MLP\_Rollercoaster*, *MLP\_Ensemble2*, and *MLP\_Ensemble4*, all with an AUC of 0.95. *MLP\_Ensemble5* also had the highest MCC value in the multiclass classification task, confirming its superior performance.

To further analyze model performance, we selected four FPR levels (1%, 5%, 10%, 20%) and calculated the corresponding TPR. The consolidated results are shown in Figure 17, with detailed results for each FPR case available in the appendix.

At each FPR level, *MLP\_Ensemble5* outperformed the other models, achieving TPRs of 43.4%, 74.3%, 86.6%, and 95.8% at the 1%, 5%, 10%, and 20% FPR levels, respectively. Excluding *MLP\_Ensemble5*, the best-performing model varied by FPR level: *MLP\_WideUltra*

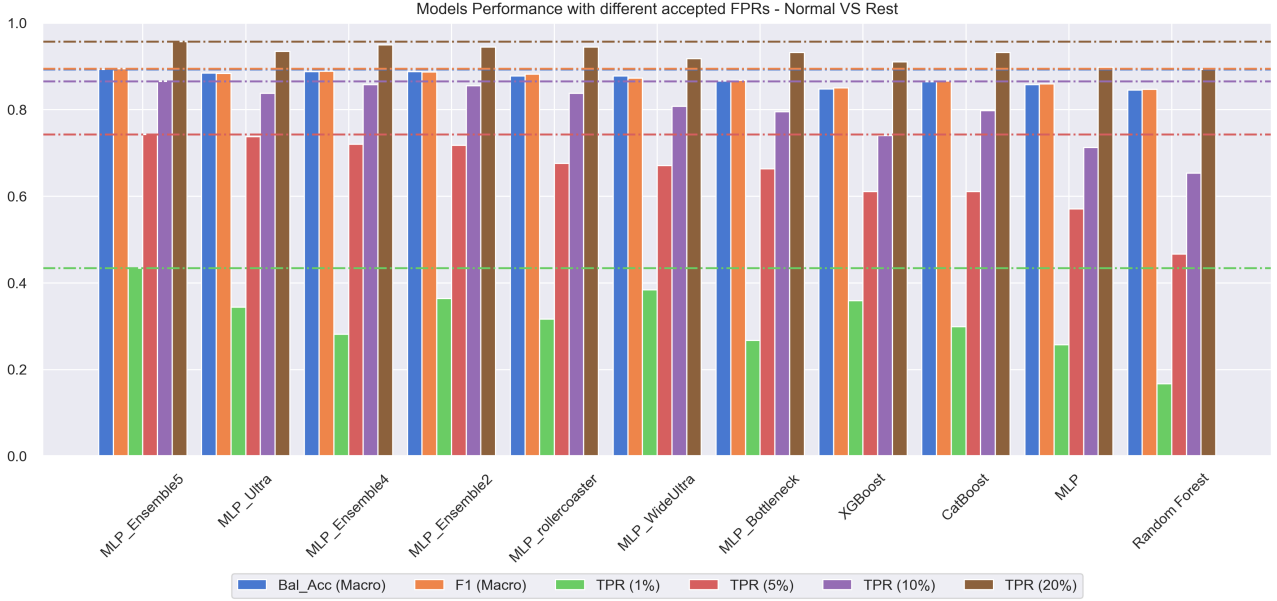


Fig. 17: TPR at different FPR levels for all models.

at 1%, *MLP\_Ultra* at 5%, *MLP\_Ensemble2* at 10%, and *MLP\_Ensemble4* at 20%.

These outcomes highlight the task's challenges in creating a model that performs well across all FPR levels and demonstrate the efficacy of a well-built ensemble model, which leverages the strengths of different models to achieve optimal performance.

### 8.3.3. Best Model Analysis

To further investigate the performance of the ensemble model, we compared the confusion matrices of the individual models with the ensemble one (Figure 18).

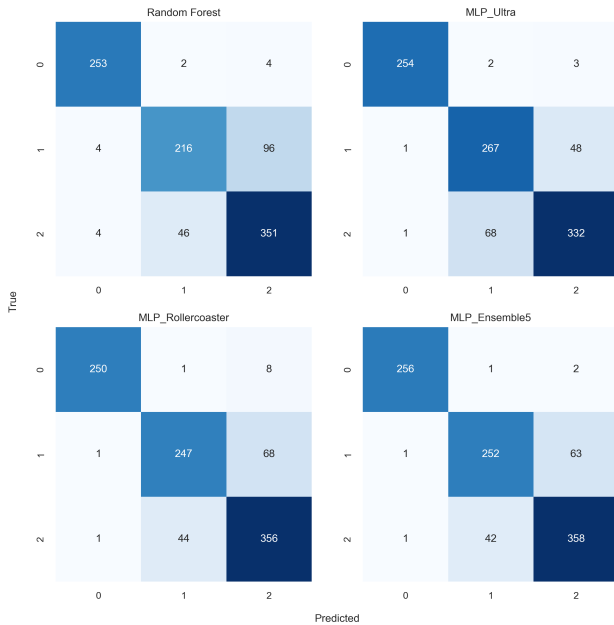


Fig. 18: Confusion matrices of the individual models and the ensemble model.

In the confusion matrix the class 2 represents the normal heartbeats, the class 1 represents the abnormal heartbeats,

and the class 0 represents the artifacts. We can see that the Random Forest and *MLP\_Ultra* are "complementary" models, with the former having higher true positive rates for the normal class and the latter a smaller false positive rate. The ensemble model combines these strengths. The contribution of the *MLP\_Rollercoaster* model is less evident, but it was experimentally shown that it contributes to the ensemble's performance. The reason may be related to the fact that the some of the samples that are misclassified by the other models are correctly classified by the *MLP\_Rollercoaster* model.

Interestingly, the *MLP\_Ultra* model classifies less abnormal heartbeats as normal than the *MLP\_Ensemble5* model. However this is due to the fact that this latter simply classifies less samples as normal. Indeed we could minimize the FPR of the normal class just by classifying all the samples as abnormal, but this would result in a very low TPR for the normal class. This evidence the importance of analyzing FPR and TPR together. In conclusion, the ensemble model is the best model for the task of classifying normal heartbeats, abnormal heartbeats, and artifacts, in a FPR/TPR trade-off maximization scenario.

## 8.4. Diseases Only Model

### 8.4.1. Results

## 9. ANALYSIS OF THE BEST MODEL

## 10. CONCLUSION

## 11. FUTURE WORKS

## 12. APPENDIX

## REFERENCES

- [1] en. URL: <https://www.kaggle.com/datasets/mersico/dangerous-heartbeat-dataset-dhd>.
- [2] P. Bentley et al. *The PASCAL Classifying Heart Sounds Challenge 2011 (CHSC2011) Results*. URL: <http://www.peterjbentley.com/heartchallenge/index.html>.
- [3] Wei Chen et al. “Deep Learning Methods for Heart Sounds Classification: A Systematic Review”. In: *Entropy* 23.6 (May 2021), p. 667. ISSN: 1099-4300. DOI: [10.3390/e23060667](https://doi.org/10.3390/e23060667).
- [4] Ali Raza et al. “Heartbeat Sound Signal Classification Using Deep Learning”. en. In: *Sensors* 19.2121 (Jan. 2019), p. 4819. ISSN: 1424-8220. DOI: [10.3390/s19214819](https://doi.org/10.3390/s19214819).