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

Ligari D. • Alberti A. [1]

[1] *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 28, 2024

**Abstract** — Heart disease remains one of the leading causes of mortality worldwide, making early diagnosis crucial. This study aims to predict heart diseases by analyzing heartbeat audio signals using machine learning and network analysis. We utilized a dataset from the PASCAL Classifying Heart Sounds Challenge 2011, which includes normal heart sounds, murmurs, extra heart sounds, extra systoles, and artifacts. Various preprocessing techniques such as noise reduction, resampling, and segmentation were applied to ensure data quality. Features were extracted using methods like Mel-Frequency Cepstral Coefficients (MFCC), Chroma, RMS, ZCR, and spectral features. Multiple machine learning models including LightGBM, XGBoost, CatBoost, Random Forest, and Multilayer Perceptron were trained and evaluated. The best performing model achieved high accuracy in distinguishing between different heart sound categories. This research highlights the potential of machine learning in cardiac diagnostics and provides a foundation for future advancements in the field.

**Keywords** ——TO BE DEFINED—

## CONTENTS

- Scikit-learn - Numpy - Pandas - Matplotlib - Seaborn - Scipy - XGBoost - CatBoost - PyTorch - Torchaudio - Librosa - TensorFlow - Keras - Shap - Imblearn - Other Utility Libraries (e.g. joblib, os, sys, etc.)

# 6. APPENDIX