# PPG-DaLiA Dataset Analysis

# 1.

# **Dataset**

Introducing the data

# Dataset

## Description

PPG-DaLiA is a publicly available dataset for PPG-based heart rate estimation. This dataset features physiological and motion data, recorded from both a wrist- and a chest-worn device, of 15 subjects while performing a wide range of activities under close to real-life conditions.

## Objective

From this data we were asked to predict the activity the patient is undergoing.

# Dataset format

○ The data of each patient is recorded in a .pkl file
○ This file contains multiple information :
- attributes of the patient (height, weight, age etc…)
- the activity executed
- signals from different captors
- indexes of the identified R-peaks
- heart rate ground truth extracted from the ECG-signal

However this data is not expressed uniformly, these features are stored in different frequencies (4 Hz, 32 Hz, 64 Hz …) thus the need to process the data.
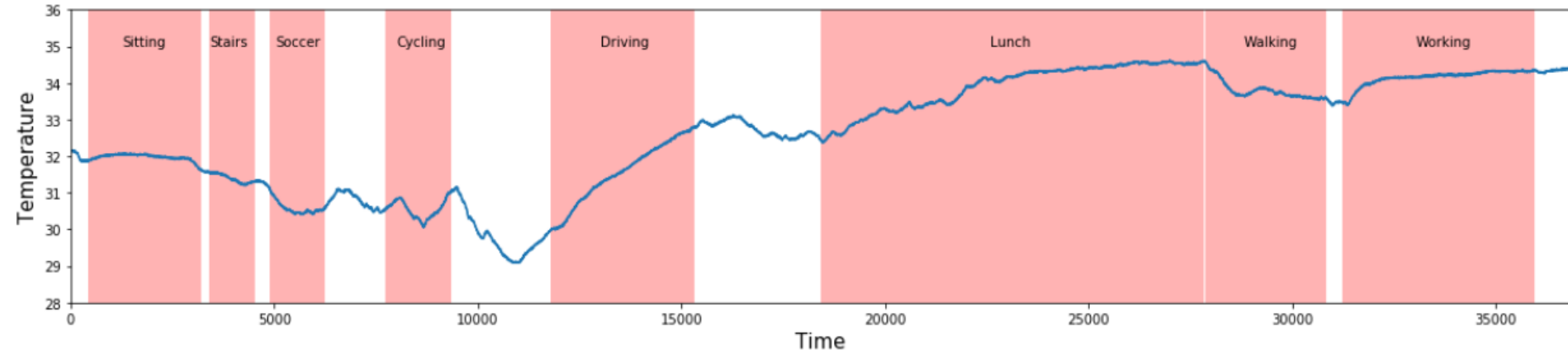
# 2.

# **Study on one patient**

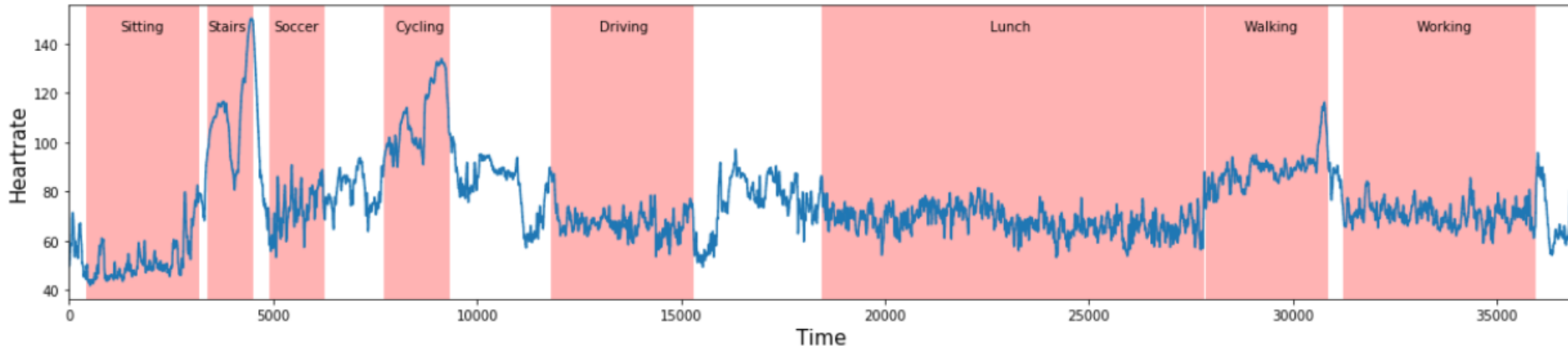Pre-processing, Visualization & Prediction

# Pre-processing

○ As mentioned before the features are stored in different frequencies. The first step was to aggregate the data and put them in the same frequency. I chose to put it all to 4Hz as it is the frequency used to record the activity (our target variable).

○ No feature engineering was needed except for the R-Peaks. Since the indexes were given, I chose to count the number of R-Peaks that occurred during our time period (4Hz so 0.25 seconds).

# Visualization



**Plot of the temperature during the different activities**

# Visualization



**Plot of the heart-rate during the different activities**
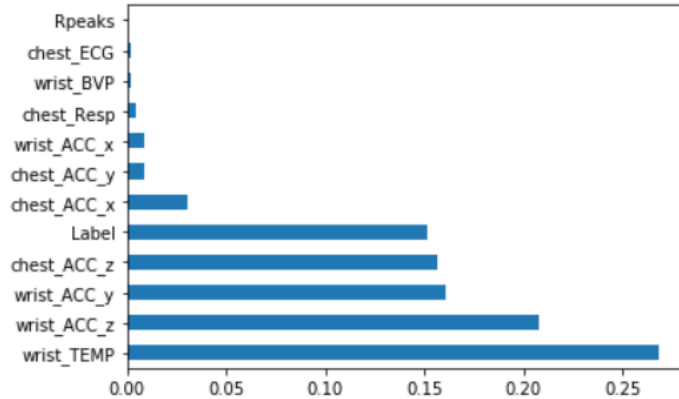
# Prediction

○ We then tried to build a model to predict the activities of the same patient. We didn't use the attributes of our patient since we are working only on one subject.

○ We've essentially testes 2 models : Decision Tree & Random Forest. Since the results of those models were great, I was more interested in seeing how they would perform when faced with different data.

# Prediction

## Decision Tree

Score : 0.975



**Feature importance**

# Prediction

## Random Forest

Score : 0.978



**Feature importance**

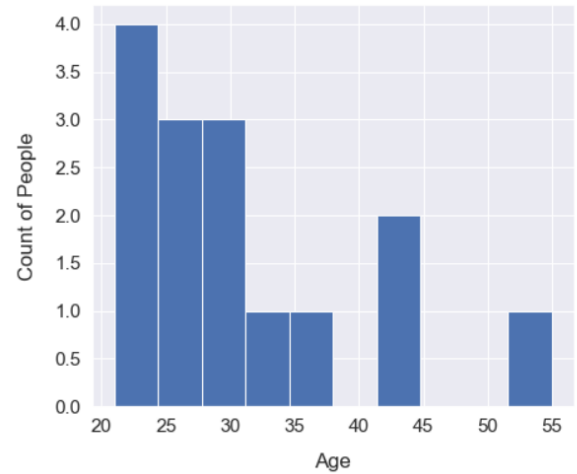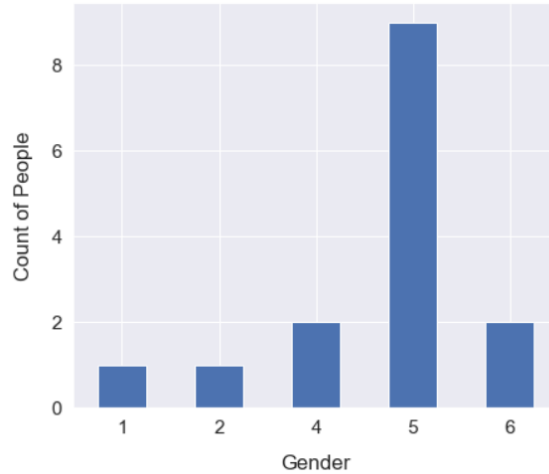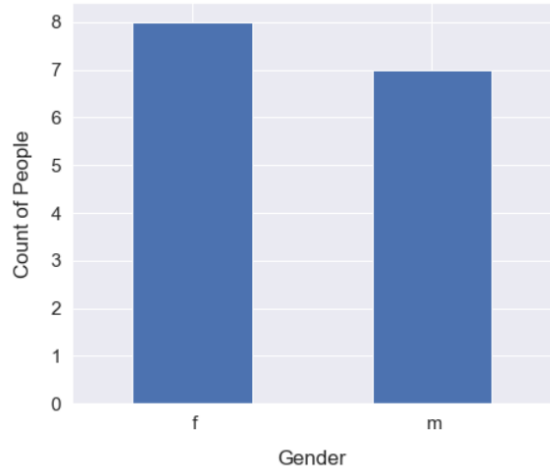# 3.

# **Study on all patients**

Pre-processing, Visualization & Prediction

# Pre-processing

o  The same steps of preprocessing will be done for each of the subjects. We will then aggregate the data of all subjects into one dataframe.
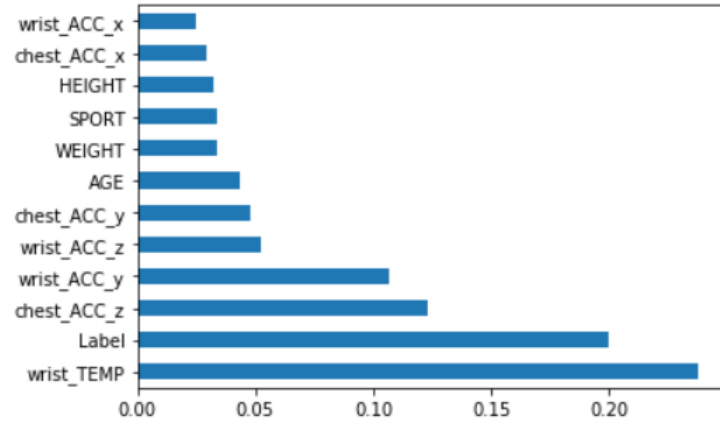
# Visualization

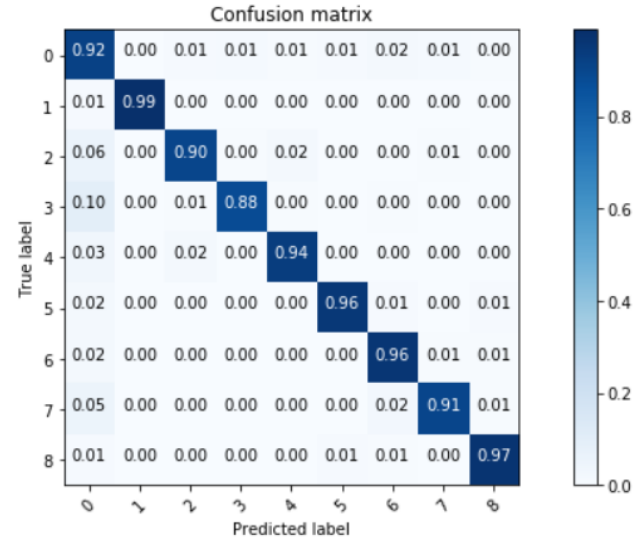*Plotting the attributes of the subjects.*
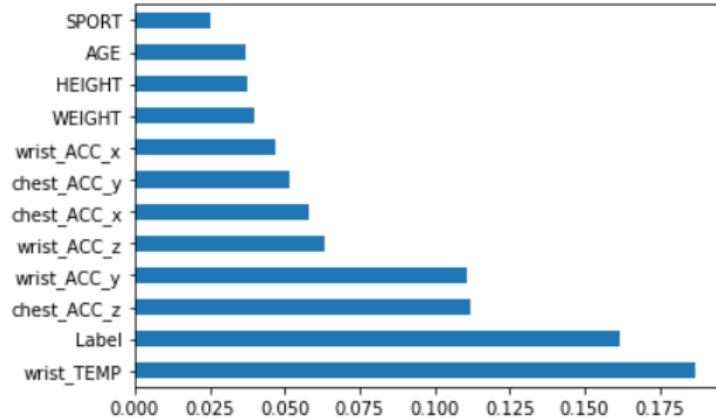
# Prediction

## Decision Tree
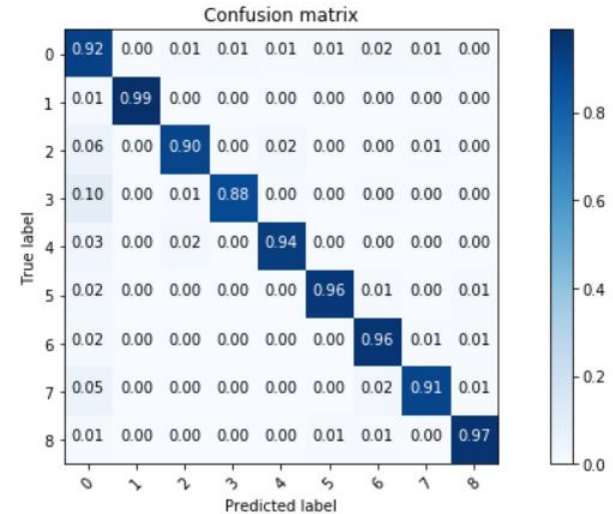
Score : 0.94



**Feature importance**

# Prediction

## Random Forest

Score : 0.962



**Feature importance**

# Prediction

○ Our model are still doing very well. My guess is that those results come from the fact that since the activities are done during a certain period, the "almost" same rows appear both in the train set and the test set. The model probably only predicts things it has already seen which means it is very likely our model is overfitting.

○ To challenge that guess, I will train the model on only 13 subjects, leaving the 2 other subjects to score our model.

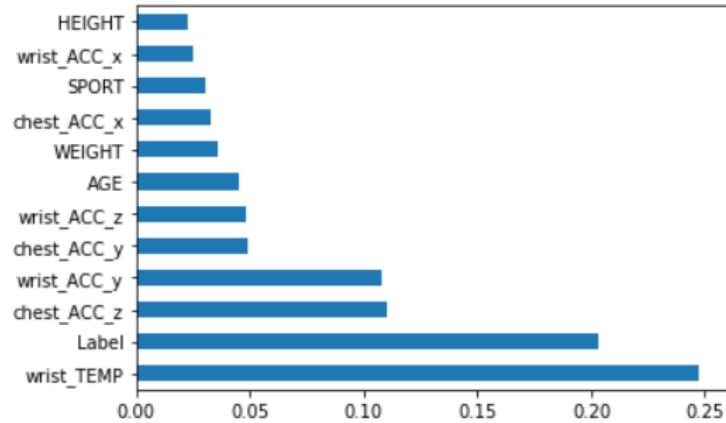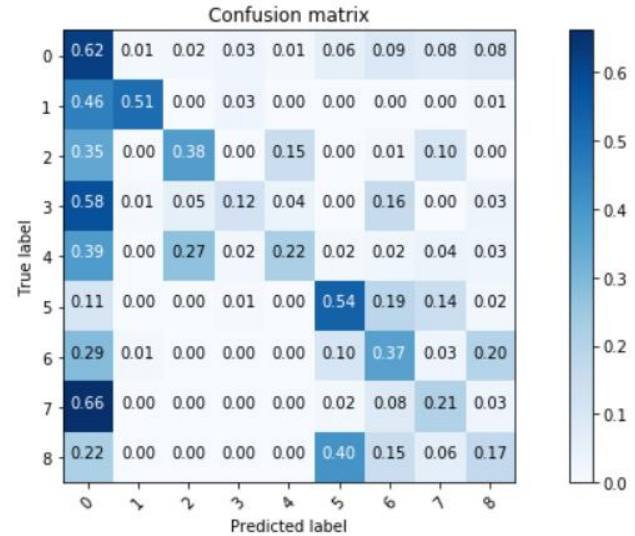# 4.

# Study on 13/2 patients
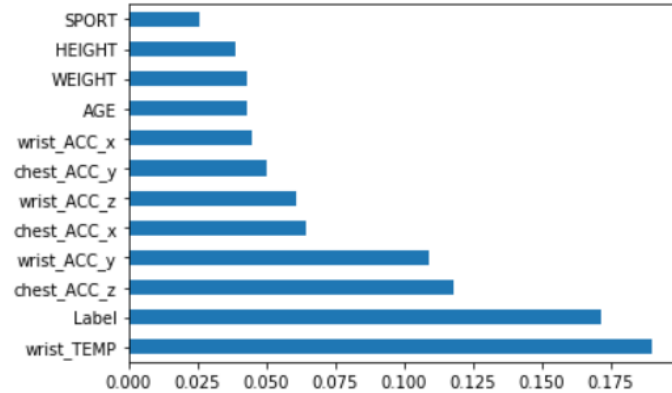
Prediction

# Prediction

## Decision Tree

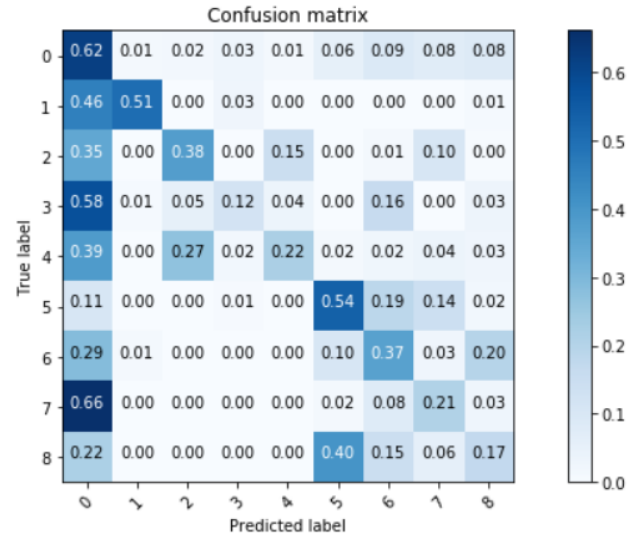Score : 0.407



**Feature importance**

# Prediction

## Random Forest

Score : 0.552



**Feature importance**

# Prediction

○   We can clearly see now that the model trained does very poorly when we face it to new subjects. The previous models were overfitting because they were only trained on 15 sets for each activities.

○   To have better results we would need data from more subjects.

END

## Credits

Special thanks to all the people who made and released these awesome resources for free:

◎ Presentation template by SlidesCarnival
◎ Photographs by Unsplash