

Capstone Proposal

January 11, 2017

Trenton Potgieter

1 A Three-Way Consensus Pipeline for Stress Level Detection

1.1 Background

As one gets older, an increasingly difficult awareness of our parent's mortality becomes a serious concern. Personally, my parents are both in their early 70's and according to a study ¹ done in 2015 by the **American Heart Association**, around 370,000 people die of heart attacks each year and is the **No. 1** cause of in the United States. In 2014, around 356,500 people experienced heart attacks out of the hospital. Of that amount only 12% survived due to emergency medical services intervention. Personally, I would not like my parents to be one the 88% who suffered from a fatal heart attack and didn't survive due to the fact that there was no intervention by emergency medical services. According to the study, there is a prevalence of almost *third* of the population at risk of *Heart Disease* leading to a *Heart Attack* as one approaches 80+ years of age. Having no personal experience in the Coronary Field of Medical research, it would be difficult for me to diagnose any potential warning signs, but with the advent of wearable technology, the mechanisms are in place to potentially aid in this early warning and detection of heart attacks. The majority of wearable technology today has the built-in ability to monitor heart rates. Therefore in this project, I proposed that this information can be uploaded or sent to a **data ingestion pipeline** that is capable of interpreting, analyzing and detecting the patterns that could be classified as symptoms of a heart attack.

Additionally, since one of the potential symptoms is the increase in heart rates. There are a number of potential factors that influence the increase in heart rate, but there are well published guidelines ² that can be used to determine anomalous patterns. If these anomalies occur, the **data ingestion pipeline** could proactively determine if a heart attack is about to *or* has occurred and alert the appropriate emergency medical response. Thus proactively preventing a fatal or near-fatal heart attack. As an added benefit, the **pipeline** mechanism can be used to monitor patients who are in *Cardiac Rehabilitation* ³.

¹(https://www.heart.org/idc/groups/ahamah-public/@wcm/@sop/@smd/documents/downloadable/ucm_480086.pdf)

²(http://www.heart.org/HEARTORG/HealthyLiving/PhysicalActivity/FitnessBasics/Target-Heart-Rates_UCM_434341_Article.jsp#.WHEiXbGZNE4)

³(<https://www.nhlbi.nih.gov/health/health-topics/topics/rehab>)

1.2 Problem Statement

For this Project, I propose creating a classification pipeline that ingests heart-rate signal data (from a simulated wearable monitor) and classifies whether the subject is in a stressful situation that could lead to *Cardiac Unrest*. Additionally, in order to prevent a "cry-wolf" scenario or *false-positives*, the pipeline employs a consensus mechanism where three classifiers must all agree on the classification. To accomplish this, the project is comprised of three stages:

1. **Ingesting signal data.** → Collect already filtered PPG ⁴ signal data with symbolic peaks (and other features) have been collected for a one-minute time segment. Each one-minute time segment is considered an observation labeled with the class *relax* or *stress*.
2. **Classification model training.**
3. **Classification model application on new, unseen data.** → The final classification is Implementing a Weighted Majority Rule Ensemble Classifier ⁵ based on the probability of the time segment observation belonging to either class, using the following:

$$\hat{y} = \arg \max_i \sum_{j=1}^m w_j p_{ij},$$

where w_j is the weight that can be assigned to the j^{th} classifier.

1.3 Datasets and Inputs

The dataset used for this Project was obtained as part of a *Proof of Concept (POC)* project in the **Dell IoT Solutions Lab** ⁶ in Santa Clara, California, where a PPG ⁷ Pulse sensor was used to measure Heart Rate Variability (HRV) ⁸ reading, similar to those found on current wearables like the **Fitbit Charge 2** ⁹. The scope of the original POC is simply to verify if the data can be extracted and filtered to detect peaks in the PPG signal for a one minute data segment. Four separate test subjects (between the ages of 68 and 76) were subjected to different stimuli to induce *stress* and *relaxing* scenarios. The one minute observations are stored in a `data.csv` file..

For the scope of this project however, I propose training three separate supervised machine learning models by applying the following methodology to create a pipeline:

1. Separate the input data into two separate repositories. One for the observations and one for the labeled output.
2. Apply **normalization** and/or **standardization** techniques to pre-process the data.
3. Define three separate models to evaluate the the data.
4. Apply the models and measure their performance.

1.4 Evaluation Metrics

Since the success criteria of the project is based on the overall **probability** of the time segment observation belonging to either class (stressed or relaxed), each individual model as well as the

⁴(<https://en.wikipedia.org/wiki/Photoplethysmogram>)

⁵(<http://scikit-learn.org/stable/modules/ensemble.html#weighted-average-probabilities-soft-voting>)

⁶(<https://www.dell.com/en-us/work/learn/internet-of-things-labs>)

⁷(<https://en.wikipedia.org/wiki/Photoplethysmogram>)

⁸(<http://www.myithlete.com/what-is-hrv/>)

⁹(<https://www.fitbit.com/charge2>)

overall consensus pipeline will be evaluated using a **Confusion Matrix**, with specific attention to the aspects:

1. **Precision:** → Measure the accuracy of each model as well as the overall pipeline.
2. **Sensitivity:** → Measure how thoroughness of each model as well as the overall pipeline.
3. **Specificity:** → Measure how well each model as well as the overall pipeline correctly measures the incorrectly classified results.

1.5 Overall Design

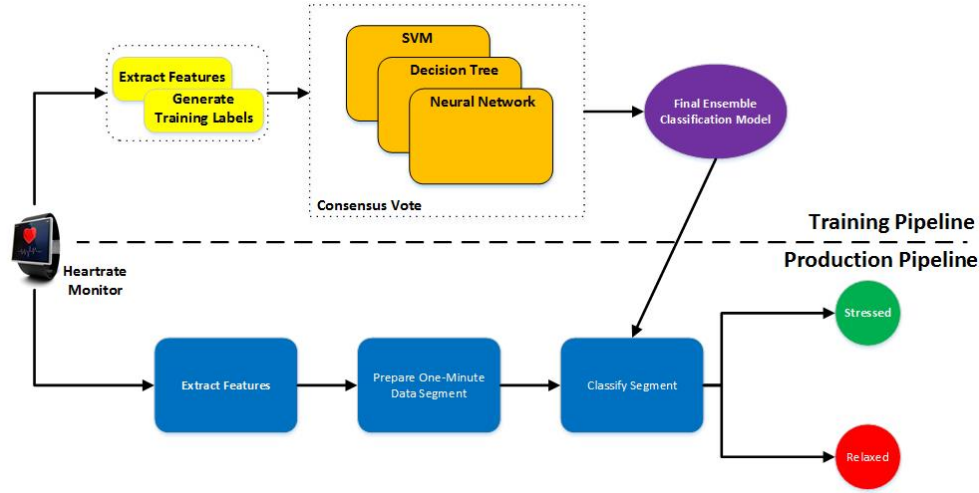


Figure 1: Training/Testing Pipeline

Figure 1 (above) provides an overview of the proposed pipeline that address the solution scope: **To determine if an individual's heart rate indicates that they are in a position of stress.**

The pipeline is separated into two specific workflows:

1. Training
2. Production

1.5.1 The Training Pipeline

The Training Pipeline is comprised of *three* specific stages.

Stage 1: Feature Extraction The first process - **Feature Extraction** - separates the incoming signal data from the heart rate monitor into two separate training data sets. The first data set are the signal observations, while the second data set are the training labels associated with each observation. The labels are further converted to a binary integer value, demarcating 1 for "relaxed" and 0 for "stressed". Additionally, in order to account for outlier variables and overfitting, the data is further standardized and scaled using the following:

Standardize:

$$X = \frac{\sum_{i=1}^n (x_i - \mu)}{\sigma}$$

Scale:

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Stage 2: Model Training Once the data has been pre-processed, three separate classifiers are trained on the data.

1. Decision Tree
2. Support Vector Machines (SVM)
3. Neural Network

Once each of these models have been trained, their resultant classification probability undergoes a consensus vote to determine the final classification probability by using the **Weighted Average Probabilities** ensemble method.

Stage 3: Final Model The last stage of the Training Pipeline is an optimized classification model that can be used for new data.

1.5.2 Testing/Production Pipeline

Like the Training Pipeline, the Production/Testing Pipeline also comprises of three stages.

Stage 1: Feature Extraction Unlike the first stage of the Training Pipeline, the data from the heart rate monitor is not separated into two data sets. Rather, the signal data is pre-processed; scaled and normalized.

Stage 2: Observation Segmentation The pre-processed data is then split into one-minute segments based on the time stamp of the data. These one-minute segments are established as a single observation of the test individual's stress level at the given time.

Stage 3: Classification The final model from the Training Pipeline is then executed against each one-minute observation segment to classify whether the test subject is stressed or relaxed.

Based on this final classification, additional future actions can be implemented that are currently outside the scope of this project.

1.6 Solution

Once created, the pipeline will be used to test and deploy the models on a sample unseen data from the test subjects and hence predict their stress levels. It is the objective of this project to re-apply the resulting pipeline to a set of new test subjects and hopefully provide a viable prototype that can preemptively warn of potential heart attacks.