# Machine Learning Prediction in mHealth Data

Team 15: Xiaoxuan Han, Hunter Ponzzebon, Gwen Eagle

## Introduction:

❖ **Dataset:**
  ➢ mHealth data (999,999 records, 14 variables);
  ➢ Time-series data on linear and angular motion on the x-,y-,z-axis during *13 different activities*
❖ **Goal:** To build a machine learning model that can accurately ascertain activity type based on multidimensional time series data
  ➢ Recurrent Neural Network (RNN)
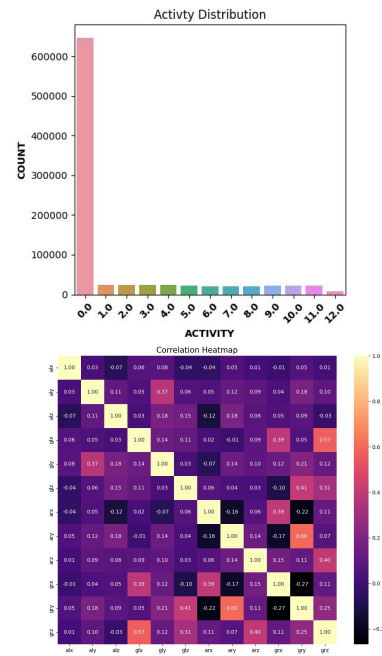  ➢ Long Short-Term Memory Neural Network (LSTMNN)

## Data Engineering & EDA :

❖ **Data Engineering:**
  ➢ Data split into 80% for training and 20% for test;
  ➢ Data normalization;
  ➢ Addressed class imbalance: undersampled *class 0* and oversampled *class 12;*
  ➢ Sequence generation
❖ **EDA:**
  ➢ No missing data;
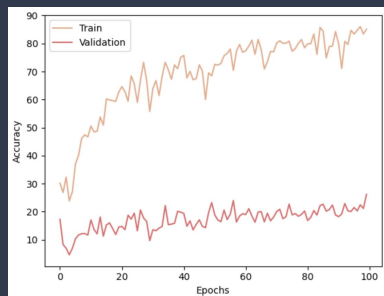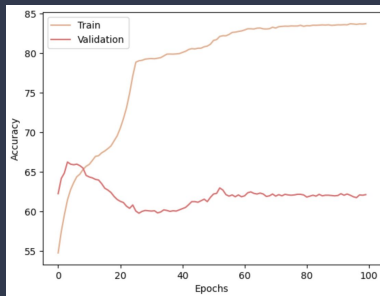  ➢ Correlation heatmap on *12 sensors*



Activty Distribution



Correlation Heatmap

# Model Comparisons

➡️

## RNN:

❖ 256 hidden units, 100 epochs, learning rate = 0.001/0.00001, batch size = 512, regularization term = 0.0001/0.00001
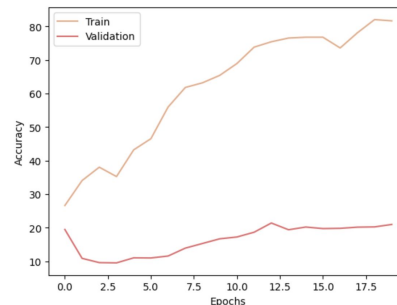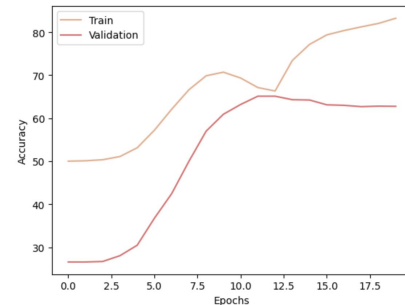


Multi-class (13 classes)



Binary class

## LSTMNN:

❖ 256 hidden units, 20 epochs, learning rate = 0.001/0.00001, a batch size = 512, regularization term = 0.0001/0.00001



Multi-class (13 classes)



Binary class

**Findings:** The RNN and LSTMNN model for the activity type performed similarly (validation accuracy of ~20%); the RNN and LSTMNN model for the presence of activity performed similarly (validation accuracy of ~62%).
**Limitations:** A high degree of overfitting might have occurred in both models, potentially due to the models trying to predict a high number of classes based on a low number of subjects, as well as a high degree of class imbalance.