# Indian Institute of Technology Kanpur



# SURGE 2023



# Research Report

**"EEG-Based Analysis of Cognitive States
and Emotional Valence: Predictive Modelling, Brainwave
Frequencies, and
Dominant Brain Locations"**

Submitted By:
**Priyansh Singh**

**SURGE Roll no. 2330339**
**Department of Chemical Engineering**
**Indian Institute of Technology Kanpur**
**Kanpur, Uttar Pradesh**

Under the guidance of:
**Dr. Tushar Sandhan**
**Assistant Professor**
**Department of Electrical Engineering**
**Indian Institute of Technology, Kanpur**
**Kanpur, Uttar Pradesh**

# CERTIFICATE

This is to certify that the project entitled "**EEG-Based Analysis of Cognitive States and Emotional Valence: Predictive Modelling, Brainwave Frequencies, and Dominant Brain Locations**" submitted by **Priyansh Singh**(2330339) as a part of Summer Undergraduate Research and Graduate Excellence 2023 offered by the Indian Institute of Technology, Kanpur, is a bonafide record of the work done by him under my guidance and supervision at the Indian Institute of Technology, Kanpur from 11th May, 2023 to 12th July, 2023.

**Dr. Tushar Sandhan**
Assistant Professor
Department of Electrical Engineering
Indian Institute of Technology Kanpur, Kanpur

# Acknowledgement

First and foremost, I am grateful to the institute IIT Kanpur, and my Project Supervisor and Mentor Prof. Tushar Sandhan for selecting me in the prestigious internship program SURGE 2023 and giving me the opportunity to work under his guidance.

I am deeply indebted to my Professor for emphasizing and encouraging me to work in the domain of EEG signals and brainwave frequencies which was very new to me and opened new doors of opportunity for me. I am thankful to him for guiding and mentoring me on projects, whenever I was stuck and for giving me a chance to explore my calibre.

I am also thankful to Ms Swathi Pratapa (Y22 MSR student) for helping me understand various pre-processing techniques used in signal processing and also guiding me through the usage of various filters in EEG-data processing.

I am extremely thankful to the Department of Electrical Engineering, IIT Kanpur for providing me with an environment that helped me develop an interest in research and further studies.

This internship would have been impossible without the constant support and inspiration of my family, who had faith in me, much more than I could have. They guided and motivated me throughout, helping me go through difficult days and inspiring me to work harder. Lastly, I would like to thank my sister, Ms. Swapnil singh (BTech Y19 electrical engineering ) for her constant support.

**Abstract**

This report presents a comprehensive investigation into analysing cognitive states and emotion valence using electroencephalogram (EEG) data. The study encompasses four distinct tracks aimed at gaining insights into the prediction of cognitive states, exploring the relationship between EEG-based brainwave frequencies and different cognitive states, examining the association between brainwaves and their dominant locations in the brain, and predicting emotion valence states using EEG data. By leveraging EEG signals, this research aims to advance our understanding of cognitive processes and emotional experiences, contributing to various applications in cognitive science and effective computing.

**Track 1**: A public data set was utilised for this. It was made up of five subjects' combined 25-hour EEG records while they were performing a low-intensity control task. The task involved utilising the "Microsoft Train Simulator" program to steer a computer-simulated train. A Butter worth high pass filter and an STFT (Blackman window) pipeline were employed for data pre-processing. The model utilized features extracted from the EEG signals captured during focused, unfocused, and drowsy states. Finally, 4 models were compared –SVM-linear, SVM-RBF, KNN, and decision trees. The model achieved the best final accuracy of 94.31 % using K-nearest neighbours.

**Track 2**: This track analysed the relationship between EEG-based brainwave frequencies and their varying presence in different cognitive states. The same dataset of track 1 was used to cater for this. Through spectral analysis of EEG data, distinct patterns of brainwave frequencies associated with different cognitive states were identified, providing insights into attention, concentration, and alertness. The "Focused" state was marked with the dominant presence of the beta and alpha waves. The "Unfocused" state had decreased presence of the beta waves with increasingly dominant alpha and theta waves. "drowsy" state had a huge presence of the theta waves.

**Track 3**: In this track, the investigation explored the association between EEG-based brainwaves and their dominant locations in the brain. The dataset of track 1 was used for this purpose. The three brainwave frequencies: alpha, beta, and theta, were plotted using various electrodes corresponding areas of the brain. The presence of alpha waves was found to be significantly high in the parietal and occipital lobes(P7,P8,O1,O2) of the brain. Similarly, the presence of the low-beta waves was found significantly high in the frontal lobes of brain(F7,F3,AF4).

**Track 4**: A different EEG dataset was utilized for emotional valence prediction. The data was collected from two people (1 male, 1 female) for 3 minutes per state -positive, neutral, and negative. Stimuli to evoke emotions were given in the form of visual movie clips marked as positive and negative. This was used to predict emotion valence as positive, negative, or neutral based on the EEG patterns. The

results demonstrated the feasibility of using EEG data for emotion valence prediction, providing valuable insights into the neural correlates of emotions.

# Contents

# List of Figures

# 1 Introduction

Neuroscience has proved that continuous exposure to mental load may affect brain dynamics. These variations are reflected in the electroencephalography (EEG) signal. Therefore, analysis of EEG signals is often utilized for cognitive load recognition. It is a non-invasive method for cognitive load detection. EEG represents the electrical activity of the brain that varies if there is any change occurs in the input conditions. Other physiological markers can only record secondary variations created by the brain, which is the fundamental response source to any environmental input. In Track 1, such variations in EEG signal under the exposure of a low-intensity control task are detected and analyzed in order to classify the overall brain response into three relevant classes i.e. focused, unfocused and drowsy state.

Detecting mental state changes in human individuals who need to remain dormant or passive while also having to maintain a continuous significant level of concentration or attention is the problem of the day. An example of such a scenario can be supervising automated processes or systems. Another example can be controlling robotic vehicles or drones or security monitoring. Yet another example can be the long-term monitoring of aircraft pilots while under the control of the autopilot. In all these cases, non-interfering supervision of processes is desired, while also alertness and a quick reaction is required of the involved individuals. With this study, we aimed to demonstrate that it is possible to identify pure mental states such as engaged and focused attention versus detached and unfocused monitoring from EEG data and tried to develop a machine learning-based system for solving such a task. Previous studies have used extra data (e.g., task engagement index or respiration data) to support EEG signals or limited the number of mental states detected to two to achieve high accuracy. We tried to detect the differentiation of engaged/focused, detached/unfocused, and drowsing mental states experimentally in a version of continuous performance tasks with 94.30 per cent (best) accuracy by using only the EEG data.

Track 2 focused on analyzing the relationship between EEG-based brainwave frequencies and their varying presence in different cognitive states . Brainwaves are divided into delta (0–4 Hz), theta (4–8 Hz), alpha (8–12 Hz), low beta (12–16 Hz), high beta (16–25 Hz), and gamma (25–50 Hz) waves, depending on their frequency. The intensity of brainwaves varies according to the state of human mental activity, and the degree of human cognitive load can be estimated based on this measurement. Theta waves are present in deeply relaxed, inward focused states with delta waves being prevalent in a sleep state. Alpha waves appear mainly in relaxed with passive attention, beta waves appear in problem solving, and gamma waves appear mainly in more complex mental functions.

The presence of alpha waves is found to be significantly high in the parietal and occipital lobes of the brain. Thus, the EEG electrodes in that locations show a dominant presence

of alpha waves. Similarly, presence of the low-beta waves is significant in the frontal lobes of the brain. These facts were verified in our Track3.

Emotion is a complicated condition that expresses human awareness and is described as a reaction to environmental stimuli. Emotions are, in general, reactions to ideas, memories, or events that occur in our environment. Russell's 2D emotion model is used most frequently. The vertical axis represents the arousal dimension (expressing the emotional intensity of the experience, ranging from low to excitement), and the horizontal axis represents the valence dimension (showing the degree of cheerfulness or joy, ranging from negative to positive). There are four categories of emotions in the arousal-valence coordinate system. The negative emotions are represented on the left side of the coordinate and the positive emotions are shown on the right. The valence axis represents both positive and negative emotions, and the arousal axis varies from inactive to active emotions(intensities). Thus, our Task 4 focused on emotional valence states prediction as positive, negative and neutral.

## 2 Dataset

### 2.1 Track-1,2,3

It was made up of five subjects' combined 25-hour EEG records while they were performing a low-intensity control task. The task involved utilising the "Microsoft Train Simulator" program to steer a computer-simulated train. In each trial, participants used the aforementioned computer simulation programme to operate a train for 35 to 55 minutes over a largely featureless course. Specifically, during the first 10 min of each experiment, the participants were engaged in focused control of the simulated train, paying close attention to the simulator's controls, and following the developments on the screen in detail. During the second 10 min of the experiments, the participants stopped following the simulator and became de-focused. The participants did not provide any control inputs during that time and stopped paying attention to the developments on the computer screen; however, they were not allowed to close their eyes or drowse. Finally, during the third 10 min of the experiments, the participants were allowed to relax freely, close their eyes and doze off, as desired. The three mental states that were looked at in this study were passive but focused attention (0–10 minutes of the job), alert but unfocused (10–20 minutes of task), and drowsy (20– 30 minutes of task). For each participant, the experiment was conducted for 7 continuous days(35-55 minutes recording on each day). The data set consisted of 14 EEG-electrodes data.

### 2.2 Track-4

The data was collected from two people (1 male, 1 female) for 3 minutes per state - positive, neutral, and negative. We used a Muse EEG headband which recorded the TP9, AF7, AF8 and TP10 EEG placements via dry electrodes. Six minutes of resting neutral data is also recorded, the stimuli used to evoke the emotions are below. 1-minute clips were shown of all the movies to both the users.

**Table 1:** Stimulus Valence Studio Year

| Stimulus | Valence | Studio | Year |
| --- | --- | --- | --- |
| Marley and Me | Neg | Twentieth Century Fox, etc. | 2008 |
| Up | Neg | Walt Disney Pictures, etc. | 2009 |
| My Girl | Neg | Imagine Entertainment, etc. | 1991 |
| La La Land | Pos | Summit Entertainment, etc. | 2016 |
| Slow Life | Pos | BioQuest Studios | 2014 |
| Funny Dogs | Pos | MashupZone | 2015 |

# 3 Pre-processing tools

## 3.1 Butter-worth filter

- A Butterworth highpass filter belongs to the class of infinite impulse response (IIR) filters, which means it has **feedback** in its design. It is characterized by its maximally flat frequency response in the pass band, meaning it **introduces minimal distortion** to the signal within the desired frequency range.

- It is designed based on the Butterworth filter's transfer function, which is determined by the order of the filter and its cutoff frequency. In our approach we have used a **5th order filter** with a cut-off frequency of **0.16Hz**. The order of the filter defines the **steepness of the roll-off**, or how quickly the filter attenuates frequencies below the cutoff frequency. A higher order filter will have a steeper roll-off but may introduce **more phase distortion**.

- The cutoff frequency of the highpass filter determines the frequency below which the signal is attenuated. Frequencies above the cutoff frequency are passed with minimal attenuation. The Butterworth highpass filter achieves a maximally flat frequency response in the **passband**, meaning there is no peaking or ripple in the passband.

- One important characteristic of the Butterworth filter is that it has a **gradual roll-off**. This means that it does not provide a sharp transition between the passband and stopband. Therefore, it may allow some lower frequency components to leak through, particularly close to the cutoff frequency.

## 3.2 K-means

- K-means clustering is a popular **unsupervised machine learning** algorithm used for clustering and grouping data points based on their similarity. The objective of the algorithm is to partition a given dataset into **K clusters**, where K is a user-defined parameter. Each cluster represents **a group of data points that are similar to each other** in terms of their features.

- The algorithm works as follows:

  Initialization: Randomly select K points from the dataset as initial **cluster centroids**.

  Assignment: Assign each data point to the nearest centroid based on a distance metric, often using **Euclidean distance**. This step creates initial clusters.

Update: Recalculate the centroid of each cluster by taking the **mean of all the data points** assigned to that cluster.

Iteration: Repeat steps 2 and 3 until convergence. Convergence occurs when the centroids no longer change significantly or when a maximum number of iterations is reached.

Final Clustering: Once convergence is achieved, the algorithm produces K clusters, where each data point belongs to the cluster with the nearest centroid.

## 3.3 STFT(Black-man window)

- **STFT** (Short-Time Fourier Transform) is a widely used time-frequency analysis technique in signal processing and audio analysis. It allows the examination of signal characteristics in **both the time and frequency** domains simultaneously by applying the Fourier transform to **short overlapping segments of the signal**.

- The STFT divides a **longer time-domain signal** into **shorter segments**, often referred to as **windows** or frames. Each windowed segment is then **Fourier transformed** to obtain its **frequency-domain representation**. By using a sliding window across the signal, the STFT provides a **time-varying** representation of the signal's frequency content.

- The windowing function plays a crucial role in the STFT as it affects the trade-off between time and frequency resolution. The **Blackman window** is a type of windowing function that helps **reduce spectral** leakage, which is a phenomenon where the **energy of a signal leaks into neighboring frequency bins** in the frequency domain. The spectral leakage can distort the frequency content of the signal and introduce unwanted artifacts.

- The Blackman window is designed to minimize spectral leakage by providing a **smooth tapering of the signal at the edges of each windowed segment.** It has a main lobe with low side lobes, which helps reduce the impact of spectral leakage.

- The formula for the Blackman window is given by:

  $$w(n) = 0.42 - 0.5 \cos\left(\frac{2\pi n}{N-1}\right) + 0.08 \cos\left(\frac{4\pi n}{N-1}\right) \quad (1)$$

  where w(n) represents the value of the window at index n, and N is the length of the window.

- By applying the Blackman window to each segment before performing the Fourier transform in the STFT, it helps to mitigate spectral leakage and **improves the frequency resolution**. This allows for a more accurate representation of the signal's frequency content, especially **for narrowband signals** or signals with sharp spectral features.

## 3.4 Gated Recurrent Unit (GRU)

- GRU (Gated Recurrent Unit) is a type of recurrent neural network (RNN) architecture that is widely used for modeling **sequential data**. It is an extension of the more traditional LSTM (Long Short-Term Memory) network and is designed to address some of the limitations of LSTM while maintaining its ability to capture long-term dependencies in sequences.

- GRU introduces the concept of **"gates"** to control the flow of information within the network. These gates **regulate the information to be remembered, forgotten, or updated** in each time step. The GRU has two main gates: the reset gate and the update gate.

- Reset Gate: The reset gate determines **which information from the previous time step should be forgotten or ignored**. It takes into account the previous hidden state and the current input to selectively reset the hidden state.

- Update Gate: The update gate controls **how much information from the current time step should be incorporated into the hidden state**. It selectively updates the hidden state by combining the previous hidden state and the current input.

- The GRU architecture enables the model to adaptively decide whether to retain or update information based on the input and previous hidden state. This gating mechanism allows for efficient training and learning of temporal dependencies in sequential data.

# 4 Classifiers

## 4.1 Support Vector Machine (SVM)-Linear

- SVM (Support Vector Machine) is a **supervised machine learning algorithm** that can be used for both classification and regression tasks. The SVM-linear classifier specifically refers to the linear variant of SVM used for classification.

- The goal of the SVM-linear classifier is to **find an optimal hyperplane that separates the data points** belonging to different classes in a **linearly separable dataset.** The hyperplane is a decision boundary that maximizes the margin, which is the distance between the hyperplane and the closest data points from each class. SVM aims to find the hyperplane that **achieves the maximum margin while minimizing classification errors.**

- The key idea behind SVM-linear classifier is to transform the original input data into a **higher-dimensional feature space using a kernel function**. In this higher-dimensional space, **the data points become more separable by a hyperplane.** However, the SVM-linear classifier uses a linear kernel, which implicitly represents the data in the original input space without explicitly mapping it to a higher-dimensional feature space. This linear kernel simplifies the computation and allows the SVM to operate efficiently, even in high-dimensional spaces.

- To find the optimal hyperplane, the SVM-linear classifier solves an optimization problem that involves minimizing a cost function, often referred to as the **hinge loss** and maximizing the margin. The hinge loss **penalizes misclassified data points**, while the margin ensures good generalization by maximizing the separation between classes.

## 4.2 Support Vector Machine (SVM)-Radial Basis Function(RBF)

- SVM-RBF (Support Vector Machine with Radial Basis Function kernel) is an extension of the SVM algorithm that allows for **non-linear classification.**The RBF kernel is defined as a similarity measure between data points in a high-dimensional feature space. It transforms the data into a space where it **becomes more separable by a hyperplane.** The RBF kernel calculates the similarity between two data points based on the **Euclidean distance** between them. It **assigns higher similarity values** to points that are **closer together** and lower values to points that are farther apart.

- In SVM-RBF, the RBF kernel replaces the dot product used in the linear SVM. The kernel function computes the similarity between pairs of data points, and these

similarities are used to define the decision boundary. By using the RBF kernel, SVM-RBF can capture non-linear relationships and **identify complex decision boundaries** that are necessary for handling more intricate classification tasks.

- The SVM-RBF algorithm aims to find the **optimal hyperplane that separates the data points** in a way that maximizes the margin and minimizes classification errors. It achieves this by **iteratively adjusting the position** and width of the decision boundary based on the support vectors, which are the data points closest to the decision boundary.

## 4.3 K-Nearest Neighbours(KNN)

- KNN (K-Nearest Neighbors) is a popular **supervised machine learning algorithm** used for both classification and regression tasks.The fundamental idea behind KNN is that **similar data points tend to belong to the same class** or have similar output values. In the case of classification, KNN determines the class label of a new data point by considering the class labels of its K nearest neighbors in the training set. The prediction is based on the **majority class among the K nearest neighbors**. In regression, KNN calculates the average or **weighted average of the output values of its K nearest neighbors** to predict the output value for the new data point.

- The K in KNN represents the number of neighbors considered for prediction. It is a hyperparameter that needs to be defined by the user. Choosing an appropriate value for K is important, as a **small K may lead to overfitting**, while a **large K may introduce more bias** into the predictions.

- The distance metric plays a crucial role in KNN for measuring the similarity between data points. Common **distance metrics** include **Euclidean distance, Manhattan distance, and Minkowski distance.** The choice of distance metric depends on the nature of the data and the problem at hand.

## 4.4 Decision Trees

- Decision trees are **supervised machine learning algorithms** that can be used for both classification and regression tasks.The structure of a decision tree resembles an upside-down tree, where each **internal node represents a decision based on a feature**, and each **leaf node represents a predicted outcome or class label.** The decision at each internal node is made by evaluating a specific feature, and the path followed through the tree leads to the final prediction at a leaf node.

- The decision tree learning process involves **recursively splitting the data based on features** to create nodes and determine the optimal decision boundaries. The goal is to find the **best feature and the best splitting criteria** at each node, which maximizes the information gain or minimizes the impurity in the resulting subsets.

# 5 Implementation and Methodology

The entire code has been divided and presented in subsections for easy understanding. The overall problem statement had four tracks,out of which the Track 1 and Track 4 use different datasets and have prediction tasks.While the track 2 and track3 use the same dataset and preprocessing as the track 1.

## 5.1 Track 1

### 5.1.1 Data Extraction

The dataset had 25 channels out of which 14 channels numbered from 4 to 17 were the eeg channels.They were plotted and those with no signal or noisy signals were disregarded.and finally on ly the 7 channels(F7 , F3 , P7 , O1 , O2 , P8 , AF4) were used.Then data was extracted from these channels in focus,unfocus and drowsy named dictionaries having (0-10 min),(10-20 min),(20-30 min) data of the experiment. The first 2 days of each subject were ignored thinking that the subject was not used to the experiment and the data after the first two days would be more relevant. Finally each state above had 5 days data for first 3 subjects and 4 days for next 2 subjects(as the dataset had only 6 days data for the last 2 subjects while for the first 3 subjects we had 7 days data). The following days belonged to each user: Subject 1: 3,4,5,6,7(1 and 2 ignored) Subject 2:10, 11, 12, 13, 14(9 and 8 ignored) Subject 3: 17,18,19,20,21(15 and 16 ignored) Subject 4: 24,25,26,27(22 ,23 and 28 ignored) Subject 5 :31,32,33,34(29,30,35 ignored) Thus ,Overall each state had 23 files of eeg data.

```python
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import scipy.io
import matplotlib.pyplot as plt
import os
from scipy import signal
from scipy.fft import fft, fftshift



file_names=[]
for dirname, _, filenames in os.walk('/content/drive/MyDrive/EEG Data/
    EEG Data'):
    for filename in filenames:
        file_names.append(os.path.join(dirname, filename))
        #print(os.path.join(dirname, filename))

#print(file_names)
# each trial is about 54 mins
```

```
19 import numpy as np
20 from scipy.fft import fft
21 from sklearn.cluster import KMeans
22 from google.colab import drive
23 drive.mount('/content/drive')
24
25
26 #starting from Channel 3,/0,1,2,3
27
28 # I have to plot all the 14 channels first.
29 #There are 14 channels, but the experimenter modified the headset.
30 #after removing the DC, the EEG amplitude should be around (-100,100).
31 #According to the plot, as the channels other than ['F7','F3','P7','01
      ','02','P8','AF4'] , don't have useful data thus only these are
      useful channels.
32 fig, ax = plt.subplots(14,1)
33 fig.set_figwidth(20)
34 fig.set_figheight(50)
35 #fig.set_size_inches(10,10)
36 mat = scipy.io.loadmat(file_names[0])
37 data = mat['o']['data'][0,0]
38 FS = mat['o']['sampFreq'][0][0][0][0]
39 channels = ['AF3', 'F7', 'F3', 'FC5', 'T7', 'P7', '01', '02', 'P8', 'T8'
      , 'FC6', 'F4', 'F8', 'AF4']
40 for i in range(14):
41     data_ave = data[5000:15000,i+3]-np.mean(data[5000:15000,i+3])
42     ax[i].plot(data_ave)
43     ax[i].set_title(channels[i])
44     ax[i].set_ylim(-100,100)
45
46
47 #F7,F3,P7,01,02,P8,AF4
48 useful_channels=[4,5,8,9,10,11,16]
49 useful_channels_names=[ 'F7', 'F3', 'P7', '01', '02', 'P8', 'AF4']
50 fig,ax = plt.subplots(8)
51 fig.set_size_inches(20,40)
52 j=0
53 for i in useful_channels:
54     data_ave = data[5000:15000,i]-np.mean(data[5000:15000,i])
55     ax[j].plot(data_ave)
56     ax[j].set_ylim(-200,200)
57     ax[j].set_title(channels[i-3])
58     j=j+1
59
60
61 marker=128*60*10
62 #marker here is determinig the 10 minutes interval
```

```
63  #delete file #28 because it doesnot have enough data
64  useful_file_index =
        [3,4,5,6,7,10,11,12,13,14,17,18,19,20,21,24,25,26,27,31,32,33,34]
65  #useful_file_index = np.arange(1,35)
66  chan_num=7
67  trail_names=[]
68  data_focus={}
69  data_unfocus={}
70  data_drowsy={}
71  focus={}
72  unfocus={}
73  drowsy={}
74  #for i in useful_file_index:
75  i=1
76  # here after enumerating the fimenames and getting their no in file name
        say 25 from eeg_record25.mat
77  #by splitting along the d and the .
78  # then by specifying the markers we divide the data in 0-10,10-20,20-30
        as focused ,unfocused and drowsy state
79  for index,filename in enumerate(filenames):
80      if int(filename.split('d')[1].split('.')[0]) in useful_file_index:
81          mat = scipy.io.loadmat(file_names[index])
82          trail_names.append(filename.split('.')[0])
83          data_focus[trail_names[-1]]=mat['o']['data'][0,0][0:marker,
        useful_channels].copy()
84          data_unfocus[trail_names[-1]]=mat['o']['data'][0,0][marker:2*
        marker,useful_channels].copy()
85          data_drowsy[trail_names[-1]]=mat['o']['data'][0,0][2*marker:3*
        marker,useful_channels].copy()
86          focus[trail_names[-1]]=mat['o']['data'][0,0][0:marker,
        useful_channels].copy()
87          unfocus[trail_names[-1]]=mat['o']['data'][0,0][marker:2*marker,
        useful_channels].copy()
88          drowsy[trail_names[-1]]=mat['o']['data'][0,0][2*marker:3*marker,
        useful_channels].copy()
```

### 5.1.2 Filtering

The data was then passed through a butterman highpass filter of 5th order and 0.16Hz cut-off frequency.

```
1  # build 5 order high pass filter
2  from scipy.signal import butter, lfilter, freqz
3  # ----- ----- ----- -----
4  def butter_highpass(cutoff, fs, order=5):
5      nyq = 0.5 * fs
6      normal_cutoff = cutoff / nyq
```

```
7      b, a = signal.butter(order, normal_cutoff, btype='high', analog=
       False)
8      return b, a
9
10  def butter_highpass_filter(data, cutoff, fs, order=5):
11      b, a = butter_highpass(cutoff, fs, order=order)
12      x = signal.filtfilt(b, a, data)
13      y = signal.filtfilt(b, a, x)
14      return y
15
16  # High Pass 0.16HZ
17  row, col = data_focus['eeg_record3'].shape
18  for name in trail_names:
19      for i in range(col):
20          data_focus[name][:,i]=butter_highpass_filter(data_focus[name][:,
       i], 0.16, 128, 5)
21          data_unfocus[name][:,i]=butter_highpass_filter(data_unfocus[name
       ][:,i], 0.16, 128, 5)
22          data_drowsy[name][:,i]=butter_highpass_filter(data_drowsy[name
       ][:,i], 0.16, 128, 5)
23          #print(name,data_drowsy[name][:,i].shape)
```

### 5.1.3 Feature Extraction

The following methodology was applied in the feature extraction.

**Defining Frequency Bins and Feature Names:** A frequency range from 0.5 Hz to
18 Hz is divided into 0.5 Hz bins. Feature names are created by combining the channel
names with the corresponding frequency bin using an underscore .

**Validating High-Pass Filter:** Two plots are generated to validate the high-pass filter.
The first plot shows the raw signal and the filtered signal for the 'eeg_record3' data. The
second plot shows the same for the 'eeg_record33' data.

**Calculating Short-Time Fourier Transform (STFT):** The Short-Time Fourier Trans-
form (STFT) is calculated for each EEG channel using a Blackman window of size 128.
The STFT is computed at 1-second intervals within each input EEG channel, resulting
in a set of time-varying Discrete Fourier Transform (DFT) amplitudes.

**Calculating Power and Averaging:** Power is calculated by taking the squared absolute
values of the STFT. Power values are then averaged over 4 consecutive frequency bins (0.5
Hz each) and 15-second running windows to reduce dimensionality and improve temporal
resolution.

**Constructing SVM Vectors:** The power values are transformed into SVM vectors by
reshaping each channel's data at a specific frequency into a single vector. The resulting
SVM vectors are logarithmically scaled by applying a logarithmic transformation.

**Visualizing Channel-Frequency Relationships:** Plots are generated to visualize the relationship between EEG channels and specific frequency components. The plots show how each channel changes with time for a particular frequency and how each frequency component changes with frequency for a specific channel.

```python
feature_names = []
freq_range=np.arange(0.5,18.5,0.5)
symb='_'
#useful_channels_names=['F7','F3','P7','O1','O2','P8','AF4']
for index,channel in enumerate(useful_channels_names):
    for freq in freq_range:
        feature_names.append(channel+symb+str(freq))
feature_names


#validate the high pass filter
fig,ax = plt.subplots(1)
fig.set_size_inches(20,10)
color = 'tab:red'
ax.set_xlabel('X-axis')
ax.set_ylabel('Y1-axis', color = color)
ax.plot(unfocus['eeg_record3'][0:500,0], color = color,label='raw signal
    ')
ax.tick_params(axis ='y', labelcolor = color)
ax.set_ylim(3500,4100)
# Adding Twin Axes to plot using dataset_2
ax0 = ax.twinx()
color = 'tab:green'
ax0.plot(data_unfocus['eeg_record3'][0:500,0],color=color,label='
    filtered signal')
ax0.set_ylim(-150,150)
fig.legend()
plt.show()

#validate the high pass filter
fig,ax = plt.subplots(1)
fig.set_size_inches(20,10)
color = 'tab:red'
ax.set_xlabel('X-axis')
ax.set_ylabel('Y1-axis', color = color)
ax.plot(unfocus['eeg_record33'][0:500,0], color = color,label='raw
    signal')
ax.tick_params(axis ='y', labelcolor = color)
ax.set_ylim(3800,4200)
# Adding Twin Axes to plot using dataset_2
ax0 = ax.twinx()
color = 'tab:green'
```

```
40 ax0.plot(data_unfocus['eeg_record33'][0:500,0],color=color,label='
       filtered signal')
41 ax0.set_ylim(-100,100)
42 fig.legend()
43 plt.show()
44
45 # STFT was then calculated at a time step of 1 s producing a set of time
       -varying DFT
46 # amplitudes X STFT (t,  ) at 1s intervals within each input EEG channel
       .
47 t_win = np.arange(0,128)
48 M = 128
49 window_blackman=0.42-0.5*np.cos((2*np.pi*t_win)/(M-1))+0.08*np.cos((4*np
       .pi*t_win)/(M-1))#window_blackman = signal.windows.blackmanharris
       (128)
50
51 #col is 7
52 power_focus = {}
53 for name in trail_names:
54     power_focus[name]=np.zeros([col,513,601])
55
56 power_unfocus = {}
57 for name in trail_names:
58     power_unfocus[name]=np.zeros([col,513,601])
59
60 power_drowsy = {}
61 for name in trail_names:
62     power_drowsy[name]=np.zeros([col,513,601])
63
64 #the output of the stft is 513*601,1 second data will produce 1 column
       of data,there are 601
65 for name in trail_names:
66     for i in range(col):
67         f, t,y1=scipy.signal.stft(data_focus[name][:,i],fs=128, window=
       window_blackman, nperseg=128,
68                     noverlap=0, nfft=1024, detrend=False,
       return_onesided=True, boundary='zeros',
69                     padded=True)
70         f, t,y2=scipy.signal.stft(data_unfocus[name][:,i],fs=128, window
       =window_blackman, nperseg=128,
71                     noverlap=0, nfft=1024, detrend=False,
       return_onesided=True, boundary='zeros',
72                     padded=True)
73         f, t,y3=scipy.signal.stft(data_drowsy[name][:,i],fs=128, window=
       window_blackman, nperseg=128,
74                     noverlap=0, nfft=1024, detrend=False,
       return_onesided=True, boundary='zeros',
```

```
75                              padded=True)
76          power_focus[name][i,:,:]=(np.abs(y1))**2
77          power_unfocus[name][i,:,:]=(np.abs(y2))**2
78          power_drowsy[name][i,:,:]=(np.abs(y3))**2
79
80  #validate the power value
81  fig,ax = plt.subplots(col,1)
82  fig.set_size_inches(20,50)
83  for i in range(col):
84      ax[i].pcolormesh(t, f, power_focus['eeg_record18'][i,:,:],vmin=0,
    vmax=2*np.sqrt(2), shading='gouraud')
85
86  fig,ax = plt.subplots(col,1)
87  fig.set_size_inches(20,50)
88  for i in range(col):
89      ax[i].pcolormesh(t, f, power_focus['eeg_record33'][i,:,:],vmin=0,
    vmax=2*np.sqrt(2), shading='gouraud')
90
91
92  # combine bins into 0.5HZ, and keep 0-18 HZ.
93
94  num=[]
95
96  power_focus_bin = {}
97  for name in trail_names:
98      power_focus_bin[name]=np.zeros([7,36,601])
99
100 power_unfocus_bin = {}
101 for name in trail_names:
102     power_unfocus_bin[name]=np.zeros([7,36,601])
103
104 power_drowsy_bin = {}
105 for name in trail_names:
106     power_drowsy_bin[name]=np.zeros([7,36,601])
107
108 for name in trail_names:
109     for chn in range(col):
110         j=0
111         for i in range(1,144,4):
112             power_focus_bin[name][chn,j,:]=np.average(power_focus[name][
    chn,i:i+4,:],axis=0)
113             power_unfocus_bin[name][chn,j,:]=np.average(power_unfocus[
    name][chn,i:i+4,:],axis=0)
114             power_drowsy_bin[name][chn,j,:]=np.average(power_drowsy[name
    ][chn,i:i+4,:],axis=0)
115             #print(np.average(power_drowsy[name][chn,i:i+4,:],axis=0).
    shape)
```

```python
116             #if name=='eeg_record3':
117             #     if chn==0:
118             #          num.append((f[i:i+4]))
119             #     print(j)
120             j=j+1
121
122
123 #print(num)
124 #print(len(num))
125
126 # avarage over 15 seconds running window.
127
128 power_focus_ave = {}
129 for name in trail_names:
130     power_focus_ave[name]=np.zeros([7,36,585])
131
132 power_unfocus_ave = {}
133 for name in trail_names:
134     power_unfocus_ave[name]=np.zeros([7,36,585])
135
136 power_drowsy_ave = {}
137 for name in trail_names:
138     power_drowsy_ave[name]=np.zeros([7,36,585])
139
140 for name in trail_names:
141     for chn in range(col):
142         j=0
143         for k in range(0,585):
144             power_focus_ave[name][chn,:,j]=np.average(power_focus_bin[
    name][chn,:,k:k+15],axis=1)
145             power_unfocus_ave[name][chn,:,j]=np.average(
    power_unfocus_bin[name][chn,:,k:k+15],axis=1)
146             power_drowsy_ave[name][chn,:,j]=np.average(power_drowsy_bin[
    name][chn,:,k:k+15],axis=1)
147             #print(np.average(power_drowsy_bin[name][chn,:,k:k+15],axis
    =1).shape)
148             j=j+1
149
150
151 # Turn the data into a vector
152 #[252,585]
153
154 svm_focus = {}
155 for name in trail_names:
156     svm_focus[name]=np.zeros([252,585])
157
158 svm_unfocus = {}
```

17

```
159 for name in trail_names:
160     svm_unfocus[name]=np.zeros([252,585])
161
162 svm_drowsy = {}
163 for name in trail_names:
164     svm_drowsy[name]=np.zeros([252,585])
165
166 for name in trail_names:
167     for j in range(585):
168         svm_focus[name][:,j]=power_focus_ave[name][:,:,j].reshape(1,-1)
169         svm_unfocus[name][:,j]=power_unfocus_ave[name][:,:,j].reshape
    (1,-1)
170         svm_drowsy[name][:,j]=power_drowsy_ave[name][:,:,j].reshape
    (1,-1)
171     svm_focus[name]=10*np.log(svm_focus[name])
172     svm_unfocus[name]=10*np.log(svm_unfocus[name])
173     svm_drowsy[name]=10*np.log(svm_drowsy[name])
174 # now, we get the svm vector 252*585 252 rows
175
176
177 # along the horizontal direction ,it shows one channel at specific
        frequecy changes with time.
178 # each channel has 36 rows of data,
179 # the following shows 7 channels, each with 36 slots of frequency.
180 fig,ax = plt.subplots(36,1)
181 fig.set_size_inches(10,150)
182 for i in range(36):
183     ax[i].plot(svm_focus['eeg_record33'][i,:],label='focus')
184     ax[i].plot(svm_drowsy['eeg_record33'][i,:],c='r',label='drowsy')
185     ax[i].plot(svm_unfocus['eeg_record33'][i,:],c='green',label='unfocus
    ')
186     ax[i].legend()
187     ax[i].set_title("The first channel with Frequency at: "+ str((i+1)
    *0.5) )
188
189 # along the horizontal direction ,it shows one channel at specific
        frequecy changes with time.
190 # each channel has 36 rows of data,
191 # the following shows 7 channels, each with 36 slots of frequency.
192 fig,ax = plt.subplots(36,1)
193 fig.set_size_inches(10,150)
194 for i in range(36):
195     ax[i].plot(svm_focus['eeg_record5'][i,:],label='focus')
196     ax[i].plot(svm_drowsy['eeg_record5'][i,:],c='r',label='drowsy')
197     ax[i].plot(svm_unfocus['eeg_record5'][i,:],c='green',label='unfocus'
    )
198     ax[i].legend()
```

```
199     ax[i].set_title("The first channel with Frequency at: "+ str((i+1)
        *0.5) )

200

201 # along the horizontal direction ,it shows one channel at specific
        frequecy changes with time.
202 # each channel has 36 rows of data,
203 # the following shows 7 channels, each with 36 slots of frequency.
204 fig,ax = plt.subplots(36,1)
205 fig.set_size_inches(10,150)
206 for i in range(36):
207     ax[i].plot(svm_focus['eeg_record5'][i,:],label='focus')
208     ax[i].plot(svm_drowsy['eeg_record5'][i,:],c='r',label='drowsy')
209     ax[i].plot(svm_unfocus['eeg_record5'][i,:],c='green',label='unfocus'
        )
210     ax[i].legend()
211     ax[i].set_title("The first channel with Frequency at: "+ str((i+1)
        *0.5) )

212

213 fig,ax = plt.subplots(36,1)
214 fig.set_size_inches(10,100)
215 #at a specific time, each channel changes with frequency
216 for i in range(36):
217     ax[i].plot(svm_focus['eeg_record33'][:,i],label='focus')
218     ax[i].plot(svm_drowsy['eeg_record33'][:,i],c='r',label='drowsy')
219     ax[i].plot(svm_unfocus['eeg_record33'][:,i],c='green',label='unfocus
        ')
220     ax[i].legend()
221     ax[i].set_title("The first channel with Frequency at: "+ str((i+1)
        *0.5) )
222 # the peak is at the lowest frequency.
223 #Each peak is the start of the data.

224

225 #at a specific time, each channel changes with frequency
226 for i in range(0,25,5):
227     plt.plot(svm_focus[trail_names[i]][:,1])
228 # the peak is at the lowest frequency.
229 #Each peak is the start of the data.
230 plt.title("The first channel with Frequency at: "+ str((21+1)*0.5))
231 plt.show()
```

### 5.1.4 Classifier Implementations

The following methodology is being followed in the below code: **Creating Target Labels:** Target labels are created for the three cognitive states: 'focus', 'unfocus', and 'drowsy'. Three lists are created: 'label_focus' containing 0s (representing the 'fo-

cus' state), 'label_unfocus' containing 1s (representing the 'unfocus' state), and 'label_drowsy' containing 2s (representing the 'drowsy' state).

**Data Preparation for Multiple Subjects:** EEG data and target labels are prepared for multiple subjects. Separate variables are created for each subject's data and target labels (e.g., subj1, subj2, subj3, subj4, subj5).

**Training Support Vector Machine (SVM) Model for Subject 2:** The SVM model is trained using the data and target labels from subject 2. The data is split into training and testing sets using the 'train_test_split' function. The training data is then scaled using 'StandardScaler'. The SVM model with a radial basis function (RBF) kernel is trained and evaluated using the scaled training and testing data. The training and testing accuracy scores are printed.

**Training SVM Models for All Subjects:** The SVM models are trained using the data and target labels from all subjects. The data is split into training and testing sets using 'train_test_split'. The training data is scaled using 'StandardScaler'. Two SVM models are trained and evaluated: one with a linear kernel and the other with an RBF kernel. The training and testing accuracy scores are printed for each model.

**Principal Component Analysis (PCA):** PCA is applied to the scaled training data. The variance explained by each principal component is calculated, and a heatmap is created to visualize the correlation between features in the training data.

**Training SVM Model with PCA:** The SVM model with an RBF kernel is trained using the PCA-transformed training data. The training and testing accuracy scores are printed.

**Training K-Nearest Neighbors (KNN) Model:** The KNN model is trained using the scaled training data. The training and testing accuracy scores are printed.

**Training Decision Trees Model:** The decision tree classifier is trained using the scaled training data. The training and testing accuracy scores are printed.

**Training Random Forest Classifier Model:** The random forest classifier is trained using the scaled training data. The accuracy scores are calculated for different values of the maximum depth of the decision trees in the random forest, and a plot is created to visualize the relationship between depth and accuracy.

**Training SVM Model with K-Fold Cross-Validation:** The SVM model with an RBF kernel is trained using K-fold cross-validation. The training and cross-validation accuracy scores are printed for each fold, and the average scores are calculated.

```
1  #--------0
2  label_focus = [0]*585
3  #--------1
4  label_unfocus = [1]*585
5  #-------2
6  label_drowsy = [2]*585
7
```

```python
#subject is the variable for all participants

subj1_files={'eeg_record3','eeg_record4','eeg_record5','eeg_record6','
    eeg_record7'}
subj2_files={'eeg_record10','eeg_record11','eeg_record12','eeg_record13'
    ,'eeg_record14'}
subj3_files={'eeg_record17','eeg_record18','eeg_record19','eeg_record20'
    ,'eeg_record21'}
subj4_files={'eeg_record24','eeg_record25','eeg_record26','eeg_record27'
    }
subj5_files={'eeg_record31','eeg_record32','eeg_record33','eeg_record34'
    }


# I will try to use the data from all participants to train the model
target=[]
subj=np.array([]).reshape(252,0).copy()
for name in trail_names:
    subj=np.concatenate((subj,svm_focus[name]), axis=1)
    subj=np.concatenate((subj,svm_unfocus[name]), axis=1)
    subj=np.concatenate((subj,svm_drowsy[name]), axis=1)
    target = target+label_focus+label_unfocus+label_drowsy
subj=subj.T
target = np.array(target)

# This part I only train the data for subject1
target1=[]
subj1=np.array([]).reshape(252,0).copy()
for name in subj1_files:
    subj1=np.concatenate((subj1,svm_focus[name]), axis=1)
    subj1=np.concatenate((subj1,svm_unfocus[name]), axis=1)
    subj1=np.concatenate((subj1,svm_drowsy[name]), axis=1)
    target1 = target1+label_focus+label_unfocus+label_drowsy
subj1=subj1.T
target1 = np.array(target1)

target2=[]
subj2=np.array([]).reshape(252,0).copy()
for name in subj2_files:
    subj2=np.concatenate((subj2,svm_focus[name]), axis=1)
    subj2=np.concatenate((subj2,svm_unfocus[name]), axis=1)
    subj2=np.concatenate((subj2,svm_drowsy[name]), axis=1)
    target2 = target2+label_focus+label_unfocus+label_drowsy
subj2=subj2.T
target2 = np.array(target2)
```

```python
50 target3=[]
51 subj3=np.array([]).reshape(252,0).copy()
52 for name in subj3_files:
53     subj3=np.concatenate((subj3,svm_focus[name]), axis=1)
54     subj3=np.concatenate((subj3,svm_unfocus[name]), axis=1)
55     subj3=np.concatenate((subj3,svm_drowsy[name]), axis=1)
56     target3 = target3+label_focus+label_unfocus+label_drowsy
57 subj3=subj3.T
58 target3 = np.array(target3)
59
60 target4=[]
61 subj4=np.array([]).reshape(252,0).copy()
62 for name in subj4_files:
63     subj4=np.concatenate((subj4,svm_focus[name]), axis=1)
64     subj4=np.concatenate((subj4,svm_unfocus[name]), axis=1)
65     subj4=np.concatenate((subj4,svm_drowsy[name]), axis=1)
66     target4 = target4+label_focus+label_unfocus+label_drowsy
67 subj4=subj4.T
68 target4 = np.array(target4)
69
70 target5=[]
71 subj5=np.array([]).reshape(252,0).copy()
72 for name in subj5_files:
73     subj5=np.concatenate((subj5,svm_focus[name]), axis=1)
74     subj5=np.concatenate((subj5,svm_unfocus[name]), axis=1)
75     subj5=np.concatenate((subj5,svm_drowsy[name]), axis=1)
76     target5 = target5+label_focus+label_unfocus+label_drowsy
77 subj5=subj5.T
78 target5 = np.array(target5)
79
80 print('length of the target:',len(target))
81 print('the shape of the data from the subject1:', subj.shape)
82
83 # Train the data from subject1
84 from sklearn.model_selection import train_test_split
85 from sklearn import preprocessing
86 from sklearn.svm import SVC
87
88 data_train, data_test, data_train_target, data_test_target =
       train_test_split(subj5, target5, test_size=0.8, random_state=0)
89 scaler = preprocessing.StandardScaler().fit(data_train)
90 X_train_scaled = scaler.transform(data_train)
91 X_test_scaled = scaler.transform(data_test)
92
93 svm = SVC(kernel='rbf')
94 svm.fit(X_train_scaled,data_train_target)
```

```
95  print(f'The Score for Training data with SVM Model for subject2:',svm.
        score(X_train_scaled,data_train_target))
96  print(f'Score of For Test data with SVM Model for subject2 : {svm.score(
        X_test_scaled,data_test_target)}')
97
98  # Use 60% data for test
99  from sklearn.decomposition import PCA
100 data_train, data_test, data_train_target, data_test_target =
        train_test_split(subj, target, test_size=0.6, random_state=0)
101 scaler = preprocessing.StandardScaler().fit(data_train)
102 X_train_scaled = scaler.transform(data_train)
103 X_test_scaled = scaler.transform(data_test)
104 #PCA should be used on scaled data
105 pca = PCA()
106 pca.fit(X_train_scaled)
107 X_train_pca = pca.transform(X_train_scaled)
108 X_test_pca = pca.transform(X_test_scaled)
109 var = pca.explained_variance_/pca.explained_variance_.sum()
110 var[0:20].sum()
111
112 import seaborn as sns
113
114 data_corr = pd.DataFrame(X_train_scaled).corr()
115 sns.heatmap(data_corr)
116 plt.show()
117
118 import seaborn as sns
119 #sns.heatmap(pca.components_)
120 #print(sum(pca.components_[:,0]**2))
121
122 pca_map=pd.DataFrame(pca.components_,columns=feature_names,index=np.
        arange(1,253))
123 #pca.components_.shape
124 pca_map
125
126
127 #USE SVM linear model
128 svm = SVC(kernel='linear')
129 svm.fit(X_train_scaled,data_train_target)
130 print(f'The Score for Training data with SVM Linear Model for all
        subjects:',svm.score(X_train_scaled,data_train_target))
131 print(f'Score of For Test data with SVM Linear Model for all subjects :
        {svm.score(X_test_scaled,data_test_target)}')
132
133 #Use RBF model
134 svm = SVC(kernel='rbf')
135 svm.fit(X_train_scaled,data_train_target)
```

```python
136  print(f'The Score for Training data with SVM RBF Model for all subjects:
         ',svm.score(X_train_scaled,data_train_target))
137  print(f'Score of For Test datawith SVM RBF Model for all subjects : {svm
         .score(X_test_scaled,data_test_target)}')
138
139  svm = SVC(kernel='rbf')
140  svm.fit(X_train_pca[:,0:30],data_train_target)
141  print(f'The Score for Training data with SVM RBF Model for all subjects:
         ',svm.score(X_train_pca[:,0:30],data_train_target))
142  print(f'Score of For Test data with SVM RBF Model for all subjects : {
         svm.score(X_test_pca[:,0:30],data_test_target)}')
143
144  #Try KNN Model
145  from sklearn.neighbors import KNeighborsClassifier
146  neighbor = KNeighborsClassifier(n_neighbors=3)
147  neighbor.fit(X_train_scaled,data_train_target)
148  print("the score for training with data from 5 participants:",neighbor.
         score(X_train_scaled,data_train_target))
149  print("the score for test data from 5 participants:",neighbor.score(
         X_test_scaled,data_test_target))
150
151  neighbor.fit(X_train_pca[:,0:30],data_train_target)
152  print("the score for training with data from 5 participants:",neighbor.
         score(X_train_pca[:,0:30],data_train_target))
153  print("the score for test data from 5 participants:",neighbor.score(
         X_test_pca[:,0:30],data_test_target))
154
155  # decision trees
156  from sklearn.tree import DecisionTreeClassifier
157  dt = DecisionTreeClassifier(max_depth=16)
158  dt.fit(X_train_scaled,data_train_target)
159  print(f"the score for training with data from 5 participants:{dt.score(
         X_train_scaled,data_train_target)}")
160  print("the score for test data from 5 participants:",dt.score(
         X_test_scaled,data_test_target))
161
162  from sklearn.ensemble import RandomForestClassifier
163
164  train_scores = []
165  test_scores = []
166  for depth in range(1,35):
167      dt_reg = RandomForestClassifier(max_depth=depth,random_state=0)
168      dt_reg.fit(X_train_scaled,data_train_target)
169      train_scores = train_scores +[dt_reg.score(X_train_scaled,
         data_train_target)]
170      test_scores = test_scores + [dt_reg.score(X_test_scaled,
         data_test_target)]
```

```
171
172  x = list(range(1,35))
173  plt.plot(x,train_scores,c='r',label='train')
174  plt.plot(x,test_scores,c='b',label='test')
175  plt.xlabel('depth')
176  plt.ylabel('accuracy')
177  plt.title('Depth vs. Accuracy for Random Forest Classifier')
178  plt.show()
179
180
181  #from sklearn import tree
182
183  # graphviz.Source takes a graphviz object and turns it into an image.
         Graphviz is a
184  # standard graphical library/format for rendering graphs - nodes
         connected by edges -
185  # such as our decision tree below.
186  #graphviz.Source(tree.export_graphviz(dt,feature_names=feature_names))
187  #r=tree.export_text(dt, feature_names=feature_names)
188  #print(r)
189
190  from sklearn.model_selection import KFold
191  import tensorflow as tf
192  tf.random.set_seed(0)
193  kfold = KFold(n_splits=5, shuffle=True)
194
195  data_train_target=data_train_target
196
197  fold_no = 1
198
199  score_tr=[]
200  score_cv=[]
201
202  svm = SVC(kernel='rbf')
203
204  for train,cv in kfold.split(X_train_scaled,data_train_target):
205      X_tr = X_train_scaled[train]
206      Y_tr = data_train_target[train]
207      X_cv = X_train_scaled[cv]
208      Y_cv = data_train_target[cv]
209
210      tf.random.set_seed(0)
211
212      svm.fit(X_tr,Y_tr)
213      score_tr.append(svm.score(X_tr,Y_tr))
214      score_cv.append(svm.score(X_cv,Y_cv))
215
```

```
216     print(f'Score for {fold_no} Fold Training: {score_tr[-1]:.3f}')
217     print(f'Score for {fold_no} Fold cv    : {score_cv[-1]:.3f}')
218     print('-------------------------------')
219     fold_no = fold_no + 1
220
221 print(f'Score of Average For Training: {np.mean(score_tr):.3f}')
222 print(f'Score of Average For CV.     : {np.mean(score_cv):.3f}')
```

## 5.2 Track 2,3

The code here was used after basic filtering of the track 1 . In this same piece of code
the variables of different subjects,their states,and indices of different eeg channels were
used to generate a huge set of plots .Where each plot plotted the percentages of al-
pha,beta,theta waves corresponding to all the days of particular subject,with particular
state,and specific channels. The methodology followed here was: **Defining Frequency
Bands:** Frequency bands for brainwaves (theta, alpha, and beta) are defined using their
corresponding frequency ranges.
**Computing Power Spectral Density (PSD):** The Fast Fourier Transform (FFT) is
used to compute the power spectral density (PSD) of the EEG data samples. The PSD
represents the power distribution across different frequency components.
**Calculating Brainwave Powers:** The power within each frequency band (theta, alpha,
beta) is calculated by summing the PSD values within the corresponding frequency range.
**Calculating Percentage of Brainwave Powers:** The total power is calculated as the
sum of theta, alpha, and beta powers. The percentages of theta, alpha, and beta powers
are then computed by dividing each power by the total power and multiplying by 100.
**Performing K-means Clustering:** The calculated brainwave percentages (theta, al-
pha, beta) are combined into a feature vector. K-means clustering is applied to cluster
the feature vectors into three clusters (assuming 'n-clusters=3').

```
1 import numpy as np
2 from scipy.fft import fft
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans
5
6 # Assuming you have a list of 5 preprocessed EEG data samples stored in
     'eeg_data'
7 # 'eeg_data' should be a list of 1D numpy arrays
8
9 # Define frequency bands for brainwaves
10 frequency_bands = {
11     'theta': (4, 8),      # Theta band (4 Hz - 8 Hz)
12     'alpha': (8, 13),     # Alpha band (8 Hz - 13 Hz)
13     'beta': (13, 18),     # Beta band (13 Hz - 18 Hz)
```

```
14  }
15
16  # Initialize lists to store brainwave percentages for each iteration
17  theta_percentages = []
18  alpha_percentages = []
19  beta_percentages = []
20  eeg_data_files=[]
21
22  # EEG data file names
23  for i in range(5):
24
25      eeg_data_files.append('eeg_record'+str(i+10))
26
27  # Iterate over each EEG data sample
28  for i in range(5):
29      # Compute power spectral density (PSD) using FFT
30      psd = np.abs(fft(data_drowsy[eeg_data_files[i]][:,[2,5]])) ** 2
31
32      # Calculate power for theta, alpha, and beta brainwave frequencies
33      theta_power = np.sum(psd[int(frequency_bands['theta'][0]):int(
      frequency_bands['theta'][1])])
34      alpha_power = np.sum(psd[int(frequency_bands['alpha'][0]):int(
      frequency_bands['alpha'][1])])
35      beta_power = np.sum(psd[int(frequency_bands['beta'][0]):int(
      frequency_bands['beta'][1])])
36
37      # Calculate total power
38      total_power = theta_power + alpha_power + beta_power
39
40      # Calculate percentages of theta, alpha, and beta powers
41      theta_percentage = (theta_power / total_power) * 100
42      alpha_percentage = (alpha_power / total_power) * 100
43      beta_percentage = (beta_power / total_power) * 100
44
45      # Append the percentages to the respective lists
46      theta_percentages.append(theta_percentage)
47      alpha_percentages.append(alpha_percentage)
48      beta_percentages.append(beta_percentage)
49
50  # Perform K-means clustering on the brainwave percentages
51  X = np.array(list(zip(theta_percentages, alpha_percentages,
      beta_percentages)))
52  kmeans = KMeans(n_clusters=3, random_state=42)
53  clusters = kmeans.fit_predict(X)
54
55  # Plot the brainwave frequencies for each iteration
56  bar_width = 0.2
```

```
57  index = np.arange(1, 6)
58
59  plt.bar(index - bar_width, theta_percentages, width=bar_width, label='
       Theta')
60  plt.bar(index, alpha_percentages, width=bar_width, label='Alpha')
61  plt.bar(index + bar_width, beta_percentages, width=bar_width, label='
       Beta')
62
63  plt.xlabel('Subject 2')
64  plt.ylabel('Percentage(Brainwave frequencies)')
65  plt.title('Brainwave Frequencies')
66
67  # Set y-axis labels with EEG data file names
68  plt.xticks(index, [f'day {i+10}' for i in range(5)])
69
70  plt.legend()
71  plt.show()
```

## 5.3 Track 4

The methodology followed in the track 4 was as follows:

**Data Visualization:**

The code begins by visualizing a sample of the features 'fft_0_b' through 'fft_749_b' using a line plot.

**Data Preprocessing:**

The target labels are mapped to numerical values (0, 1, 2) using a label mapping dictionary. The dataset is split into training and testing sets using the train_test_split function.

**Model Architecture:**

The model architecture is defined using the TensorFlow Keras API. The input shape is determined by the number of features in the dataset. A GRU (Gated Recurrent Unit) layer is added to capture temporal dependencies in the data. A flatten layer is used to convert the output of the GRU layer into a 1D tensor. The final output layer is a dense layer with softmax activation for multiclass classification.

**Model Training:**

The model is compiled with the Adam optimizer, sparse categorical cross-entropy loss, and accuracy metric. Early stopping is applied to monitor the validation loss and stop training if there is no improvement for a specified number of epochs. ModelCheckpoint saves the best model based on validation accuracy during training. LearningRateScheduler adjusts the learning rate during training. The model is trained on the training data with validation data specified for monitoring. Training history, including accuracy and loss, is

stored for visualization.

**Model Evaluation:**

The model's accuracy is evaluated on the test data. The predicted labels are generated using the model on the test data. Confusion matrix and classification report are generated based on the predicted labels and actual labels of the test data. The confusion matrix is visualized using a heatmap. The classification report provides precision, recall, F1-score, and support for each class.

```python
import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split

import tensorflow as tf
from tensorflow.keras.callbacks import EarlyStopping,
    LearningRateScheduler, ModelCheckpoint

from sklearn.metrics import confusion_matrix, classification_report


sample = data.loc[0, 'fft_0_b':'fft_749_b']

plt.figure(figsize=(16, 10))
plt.plot(range(len(sample)), sample)
plt.title("Features fft_0_b through fft_749_b")
plt.show()

data['label'].value_counts()

label_mapping = {'NEGATIVE': 0, 'NEUTRAL': 1, 'POSITIVE': 2}

def preprocess_inputs(df):
    df = df.copy()

    df['label'] = df['label'].replace(label_mapping)

    y = df['label'].copy()
    X = df.drop('label', axis=1).copy()

    X_train, X_test, y_train, y_test = train_test_split(X, y, train_size
    =0.7, random_state=123)

    return X_train, X_test, y_train, y_test
```

```
37
38 X_train, X_test, y_train, y_test = preprocess_inputs(data)
39
40 inputs = tf.keras.Input(shape=(X_train.shape[1],))
41
42 expand_dims = tf.expand_dims(inputs, axis=2)
43
44 gru = tf.keras.layers.GRU(1024, return_sequences=True)(expand_dims)
45
46 flatten = tf.keras.layers.Flatten()(gru)
47
48 outputs = tf.keras.layers.Dense(3, activation='softmax')(flatten)
49
50
51 model = tf.keras.Model(inputs=inputs, outputs=outputs)
52 print(model.summary())
53
54 model.compile(
55     optimizer='adam',
56     loss='sparse_categorical_crossentropy',
57     metrics=['accuracy']
58 )
59
60 es = EarlyStopping(monitor='val_loss', mode='min', verbose=1, patience
    =10, restore_best_weights=True)
61 mc = ModelCheckpoint('./best_dnn_model.h5', monitor='val_accuracy', mode
    ='max', verbose=1, save_best_only=True)
62 lr_schedule = LearningRateScheduler(lambda epoch: 0.001 * np.exp(-epoch
    / 10.))
63
64 history = model.fit(
65     X_train,
66     y_train,
67     validation_split=0.2,
68     validation_data=(X_test, y_test),
69     batch_size=32,
70     epochs=100,
71     callbacks=[es, mc, lr_schedule]
72 )
73
74 plt.plot(history.history['accuracy'])
75 plt.plot(history.history['val_accuracy'])
76 plt.title('model accuracy')
77 plt.ylabel('accuracy')
78 plt.xlabel('epoch')
79 plt.legend(['train', 'test'], loc='upper left')
80 plt.show()
```

```python
81  # summarize history for loss
82  plt.plot(history.history['loss'])
83  plt.plot(history.history['val_loss'])
84  plt.title('model loss')
85  plt.ylabel('loss')
86  plt.xlabel('epoch')
87  plt.legend(['train', 'test'], loc='upper left')
88  plt.show()
89
90  model_acc = model.evaluate(X_test, y_test, verbose=0)[1]
91  print("Test Accuracy: {:.3f}%".format(model_acc * 100))
92
93  y_pred = np.array(list(map(lambda x: np.argmax(x), model.predict(X_test)
        )))
94
95  cm = confusion_matrix(y_test, y_pred)
96  clr = classification_report(y_test, y_pred, target_names=label_mapping.
        keys())
97
98  plt.figure(figsize=(8, 8))
99  sns.heatmap(cm, annot=True, vmin=0, fmt='g', cbar=False, cmap='Greens')
100 plt.xticks(np.arange(3) + 0.5, label_mapping.keys())
101 plt.yticks(np.arange(3) + 0.5, label_mapping.keys())
102 plt.xlabel("Predicted")
103 plt.ylabel("Actual")
104 plt.title("Confusion Matrix")
105 plt.show()
106
107 print("Classification Report:\n---------------------\n", clr)
```

# 6 Results and Conclusion

## 6.1 Track 1:Cognitive state prediction



**Figure 1:** Test accuracies of SVM-linear on training and testing on each subject individually

**Figure 2:** Test accuracies of various models trained and tested on all-subject combined data
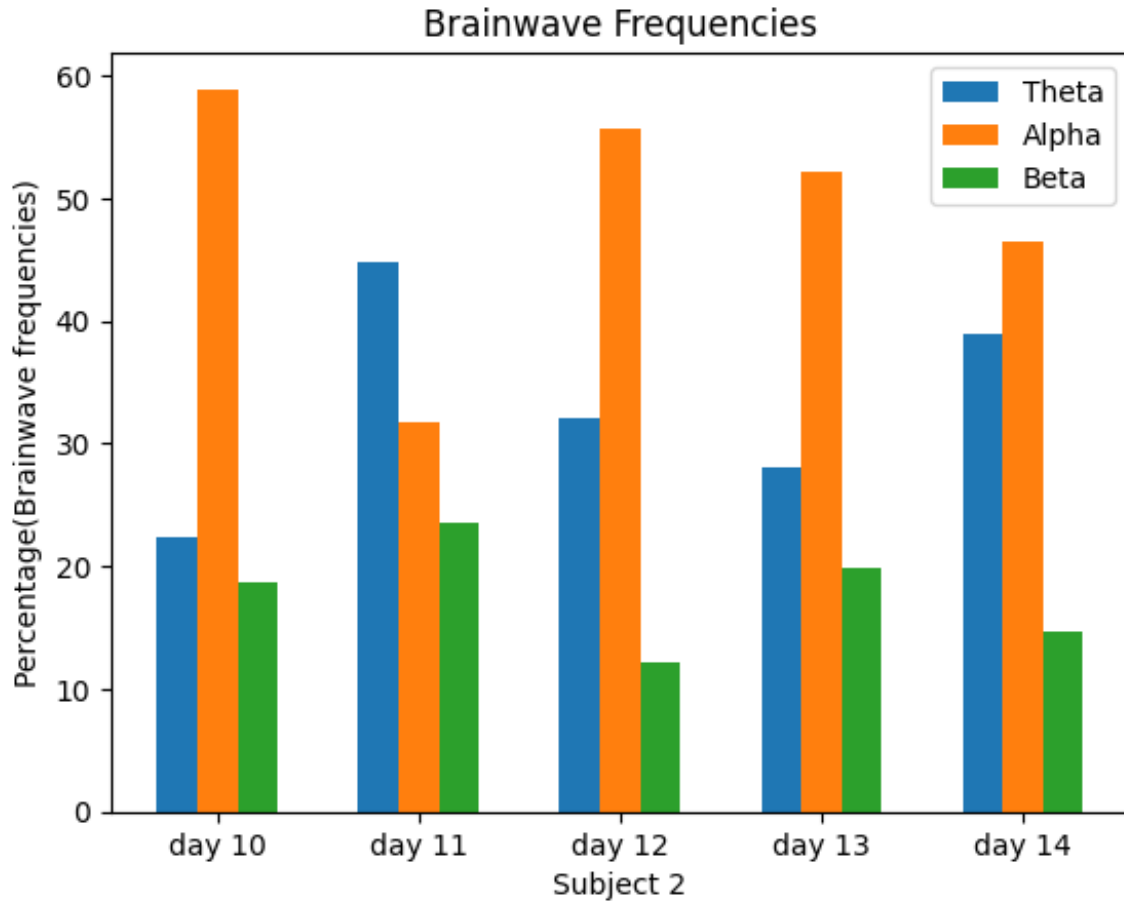
## 6.2 Track 2:Spectral presence on brainwave frequencies in different cognitive states

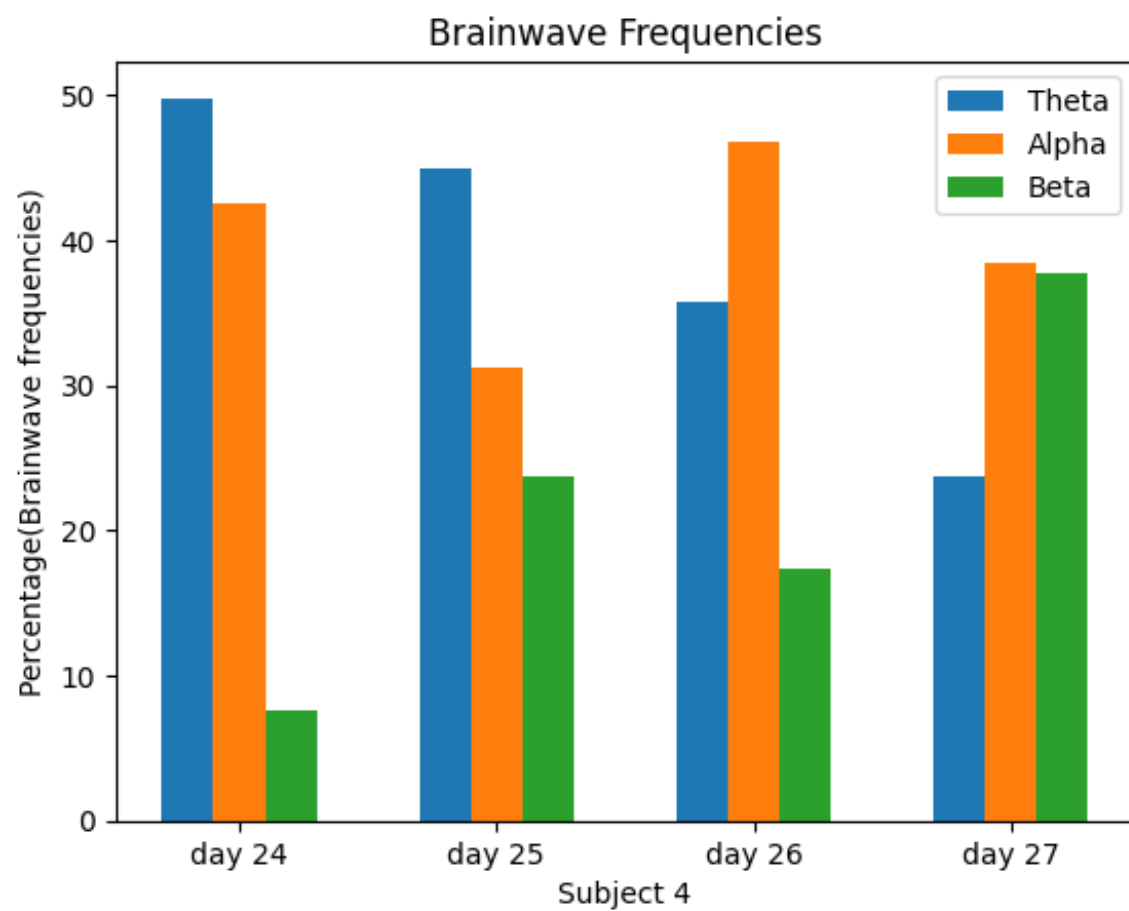**Figure 3:** Percentage of brainwave frequencies in **focussed state** of Subject 2 using all channels

**Figure 4:** Percentage of brainwave frequencies in **unfocussed state** of Subject 2 using all channels

**Figure 5:** Percentage of brainwave frequencies in **drowsy state** of Subject 2 using all channels

From the above 3 images which correspond to different states of the same subject using all channels we infer that:Drowsy state is marked by high presence of theta and alpha waves. While unfocus state show a mixture of all but slightly increased beta waves than drowsy state.while the focus state shows the highest presence of beta waves.

## 6.3 Track 3:Brainwaves and their dominant locations in brain.

### 6.3.1 alpha waves

**Figure 6:** Percentage of brainwave frequencies in drowsy state of Subject 2 using **frontal(F7,F3,AF4) channels**
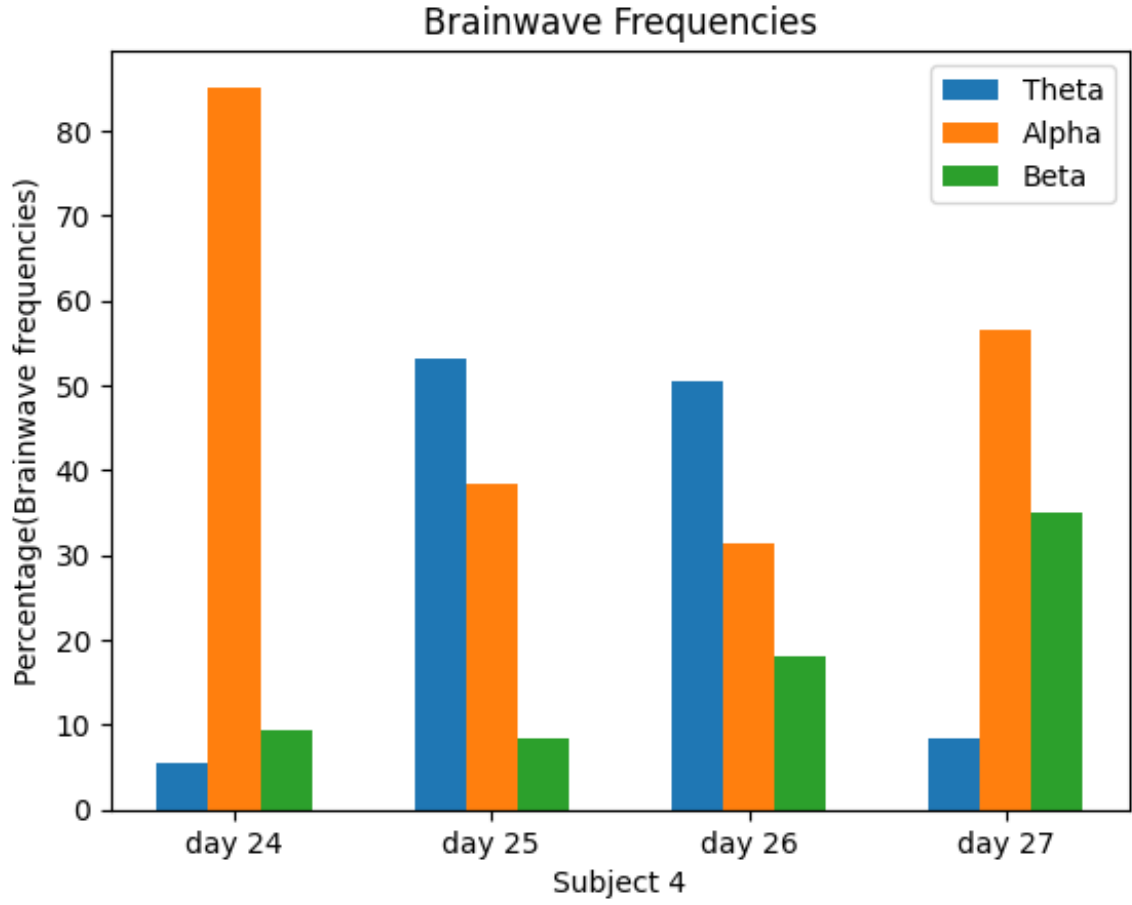
**Figure 7:** Percentage of brainwave frequencies in drowsy state of Subject 2 using **pareital and occipital(P7,P8,O1,O2) channels**

The given two plot are for drowsy states of subject 2.First plot uses only Frontal electrodes while the second one uses Pareital and occipital electrodes.We can infer that the presence of the alpha waves(orange ) is significantly high in the parietal and occipital lobes of brain.

### 6.3.2 beta waves

**Figure 8:** Percentage of brainwave frequencies in unfocus state of Subject 4 using **frontal(F7,F3,AF4) channels**
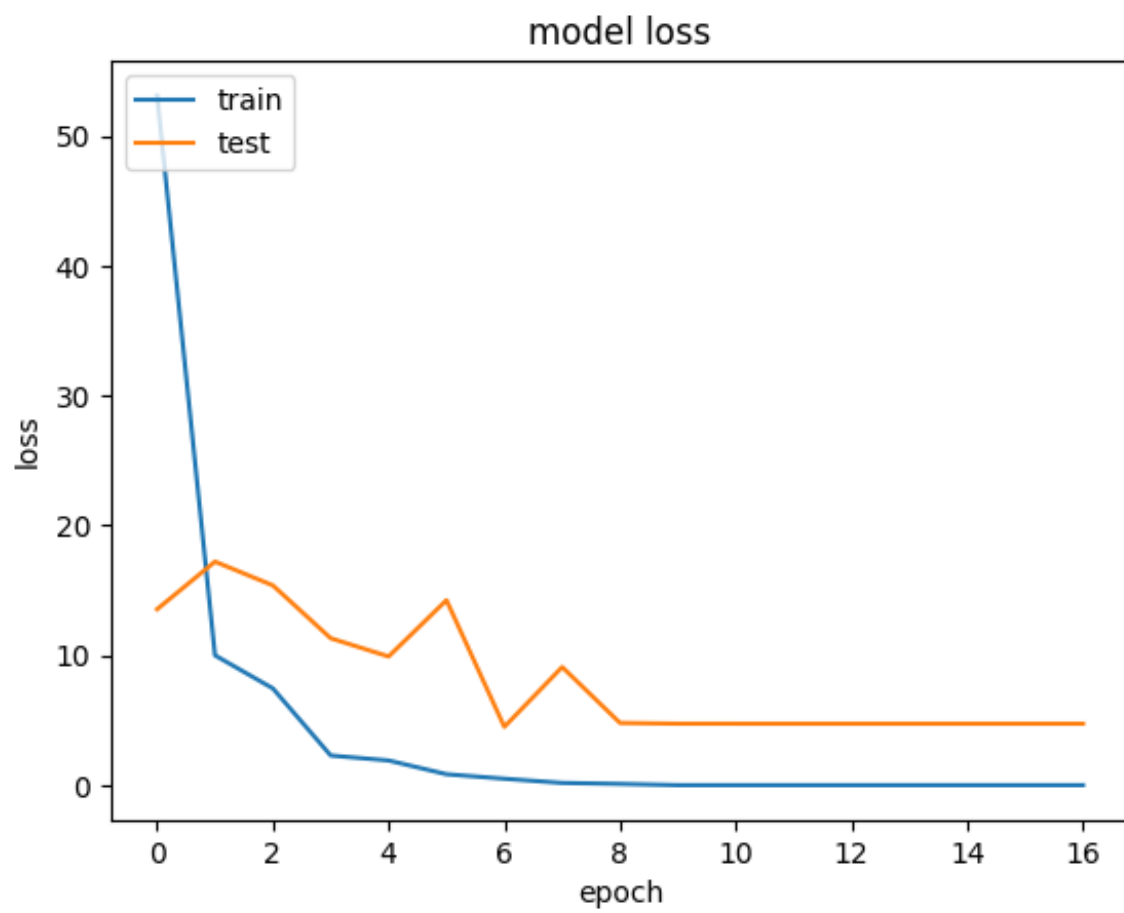
**Figure 9:** Percentage of brainwave frequencies in unfocus state of Subject 4 using **pareital and occipital(P7,P8,O1,O2) channels**
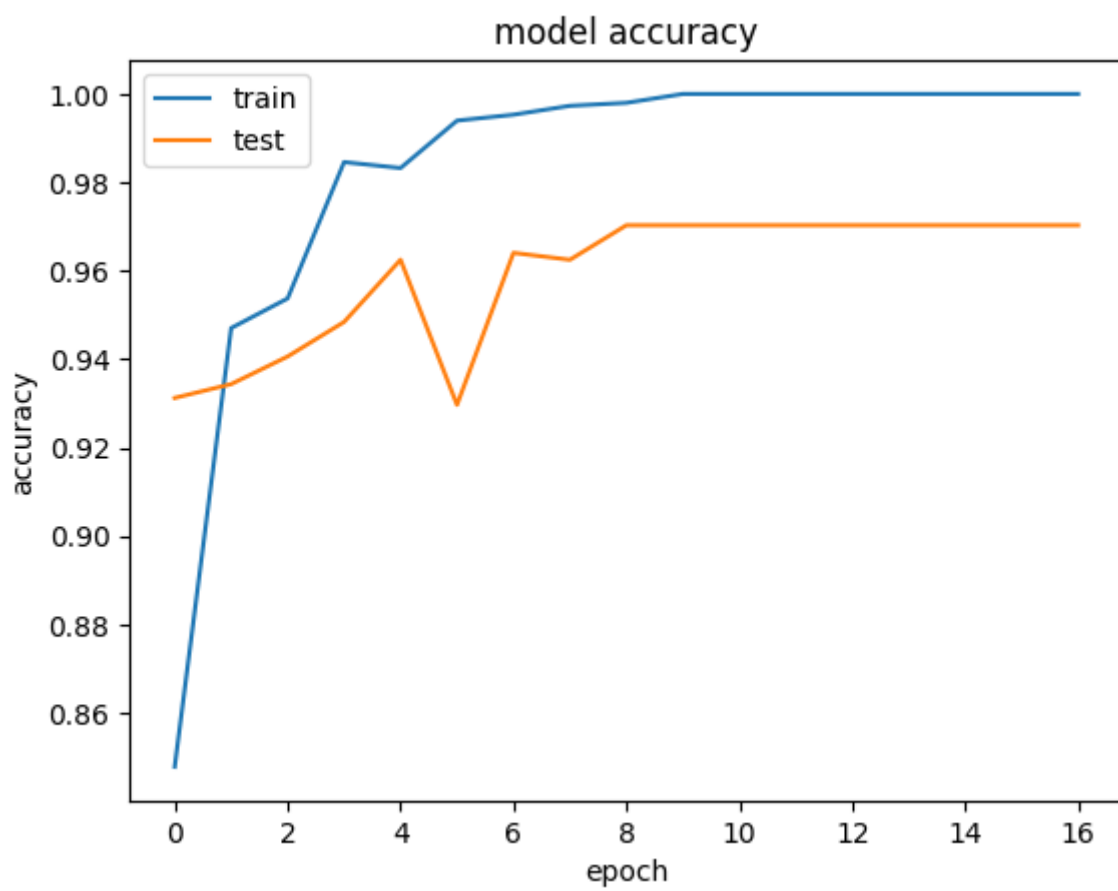
The given two plot are for unfocus states of subject 4.First plot uses only Frontal electrodes while the second one uses Pareital and occipital electrodes.We can infer that the presence of the low-beta waves(green ) is significantly high in the frontal lobes of brain.

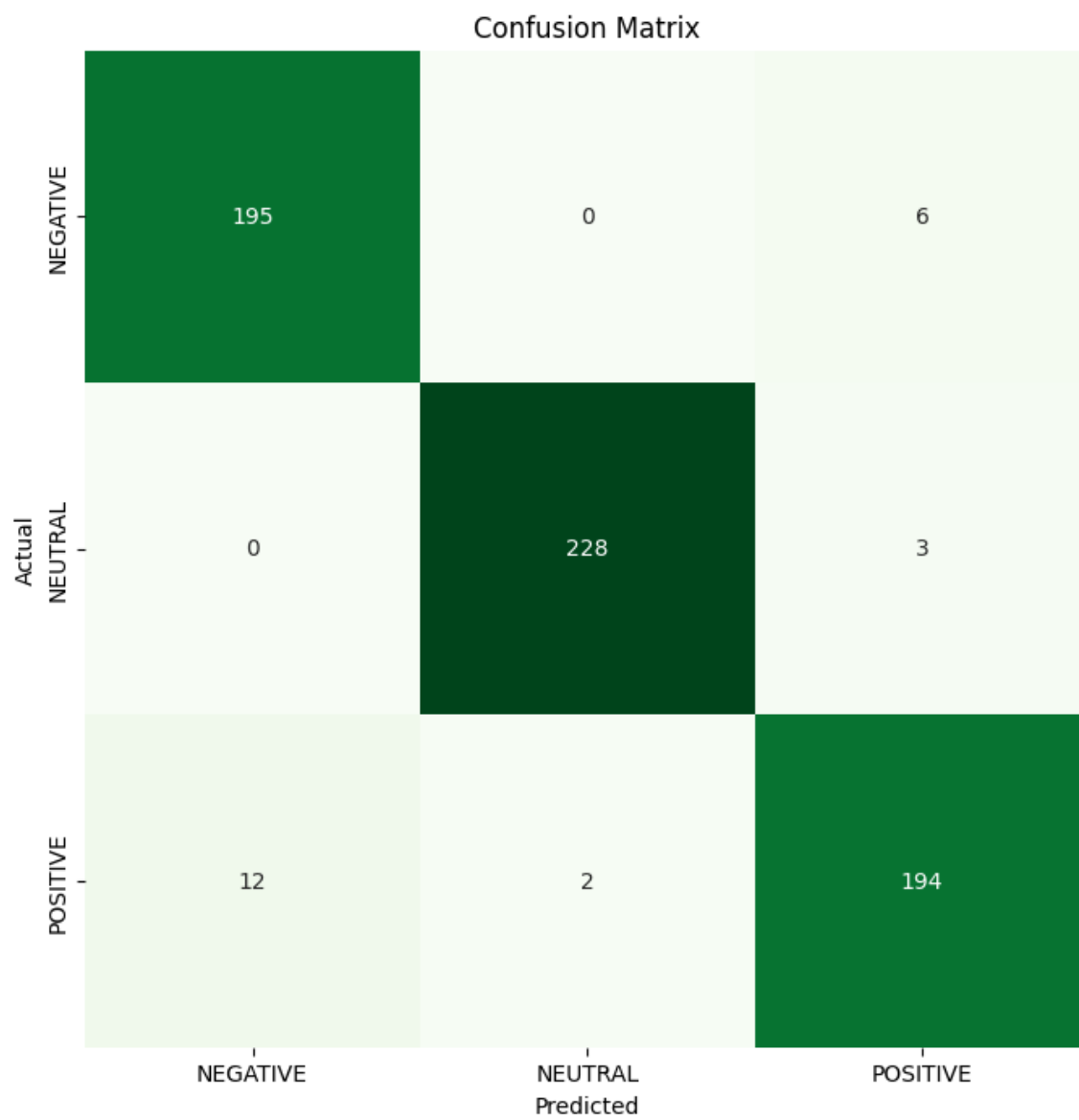## 6.4 Track 4:Emotional valence states prediction

**Figure 10:** Model loss vs epochs

**Figure 11:** Model accuracy vs epochs

**Figure 12:** Confusion matrix describing the positive ,negative ,neutral states as predicted and actual markers

# 7 References

[1] Distinguishing mental attention states of humans via an EEG- based passive BCI using machine learning methods.BY- Çigdem Inan Acı, Murat Kayaa , Yuriy Mishchenko

[2] Prediction of Cognitive Load from Electroencephalography Signals Using Long Short-Term Memory Network .BY-Gilsang Yoo, Hyeoncheol Kim and Sungdae Hong

[3] Mental arithmetic task load recognition using EEG signal and Bayesian optimized K-nearest neighbor Lakhan Dev Sharma1 • Himanshu Chhabra2 • Urvashi Chauhan3 • Ritesh Kumar Saraswat2 • Ramesh Kumar Sunkaria

[4] Mental Emotional Sentiment Classification with an EEGbased Brain-Machine Interface ,Jordan J. Bird,Anikó Ekárt,Christopher D. Buckingham,Diego R. Faria.

[5] Human emotion recognition from EEG-based brain–computer interface using machine learning: a comprehensive review Essam H. Houssein, Asmaa Hammad  Abdelmgeid A. Ali Neural Computing and Applications