

In [19]:

```
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
import warnings
import re
import os
from IPython.core.display import display, HTML
from tqdm.auto import tqdm

plt.style.use('seaborn')
%matplotlib inline
warnings.filterwarnings('ignore')
display(HTML("<style>.container { width:90% !important; }</style>"))
```

heart rate EDA

피험자 번호대로 0번부터 번호를 매긴 후 출력해본 결과,

2, 4, 30번 피험자는 [시작 시간, 끝 시간] 구간이 두 번 중복해서 나왔음. 그래서 마지막 [시작 시간, 끝 시간] 구간만 저장했음.

In [3]:

```
# heart rate 데이터 프레임 리스트
heart_rate = []

dir_name = '' # 파일 디렉토리 경로
path = '' # 데이터 파일 경로
file_list = os.listdir(path)

# 주어진 파일 이름에서 숫자 부분을 추출
# 파일 이름에서 정수 값을 추출하여 파일을 숫자순으로 정렬하는 데 사용
def get_number(filename):
    return int(re.search(r'\d+', filename).group())

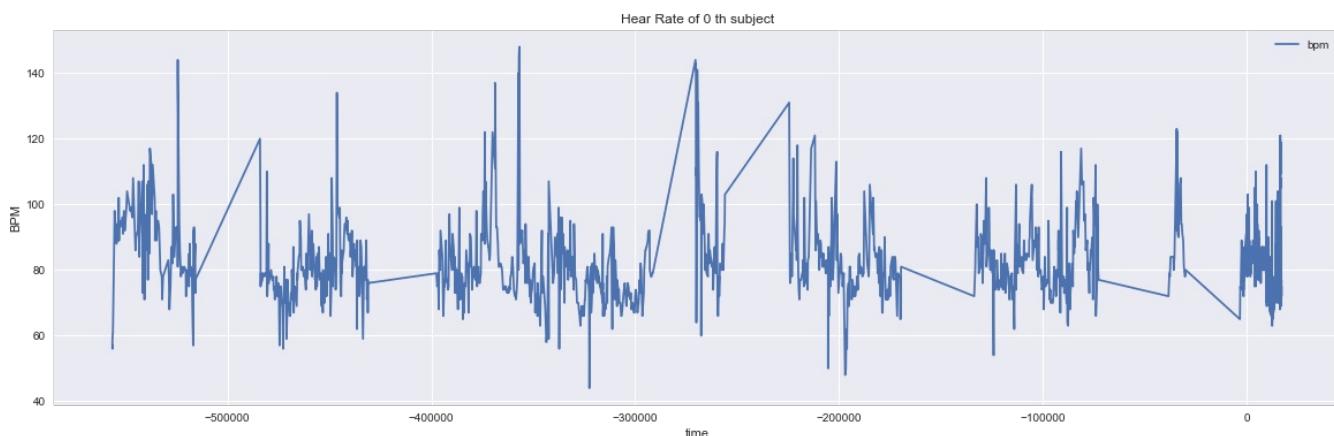
# 리스트 파일 이름 순으로 정렬
file_list_txt = sorted([file for file in file_list], key=get_number)

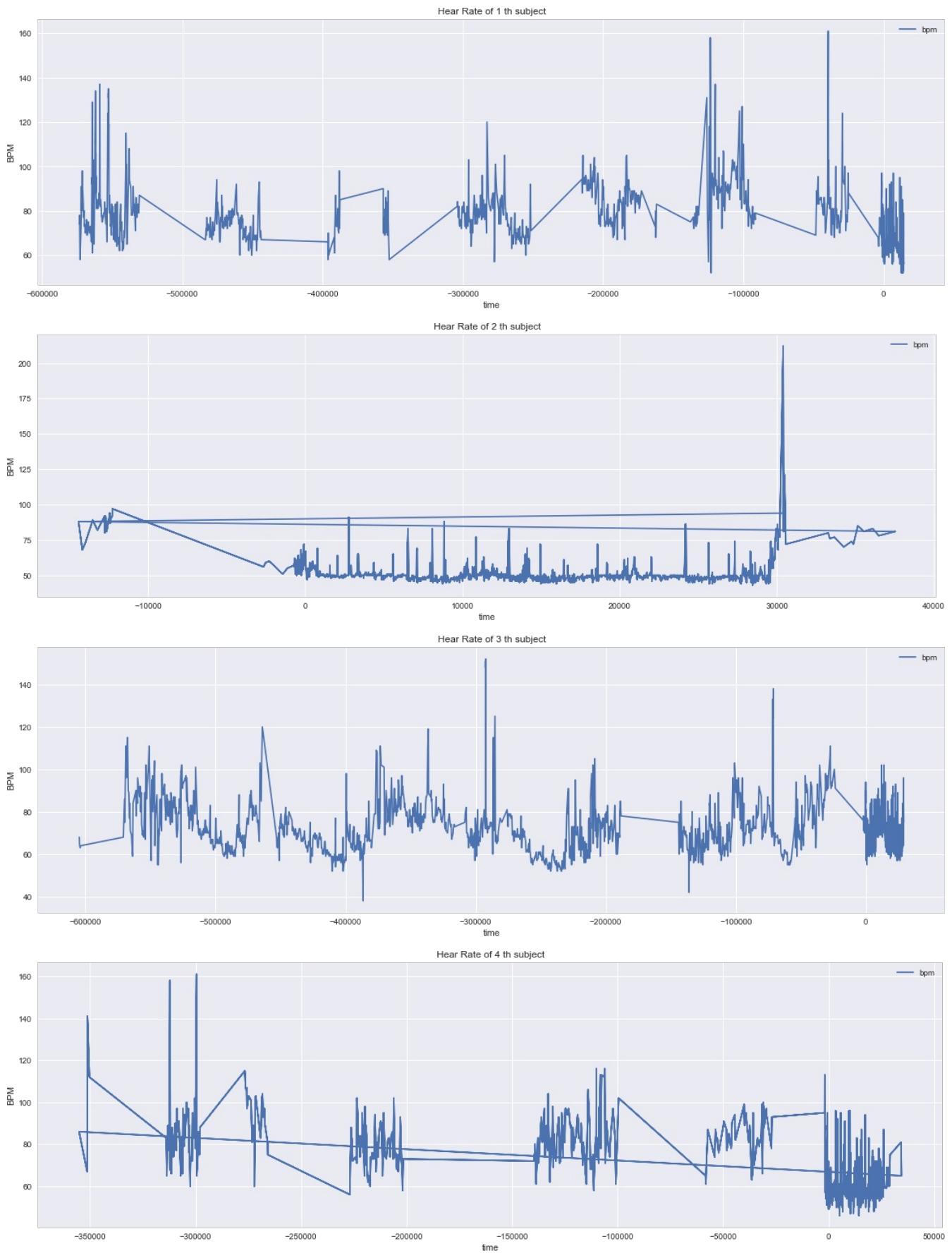
# 텍스트 파일을 csv 파일로 저장
for i in file_list_txt:
    df = pd.read_csv(dir_name+i, header=None, names=['time', 'bpm'])
    heart_rate.append(df)
```

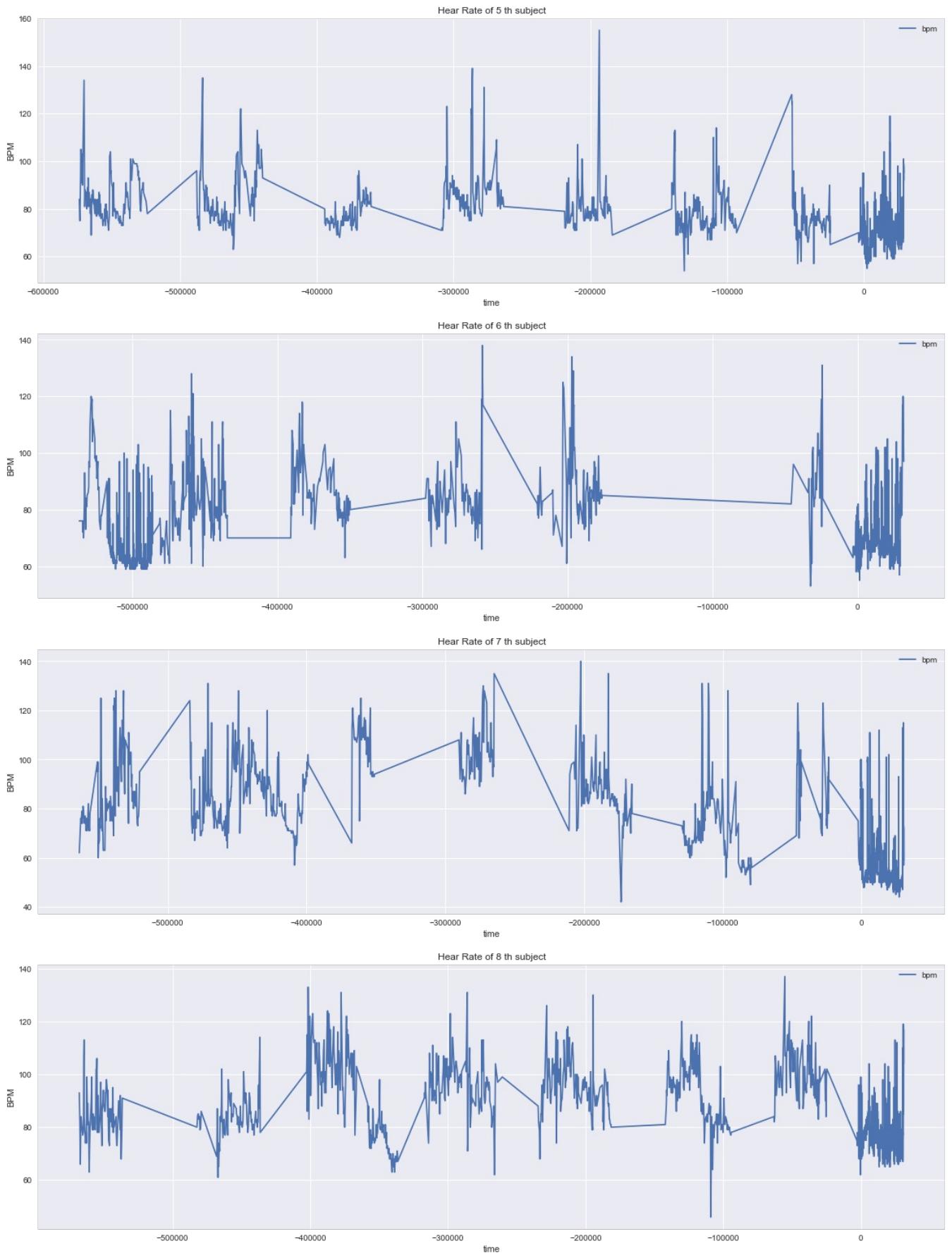
In [4]:

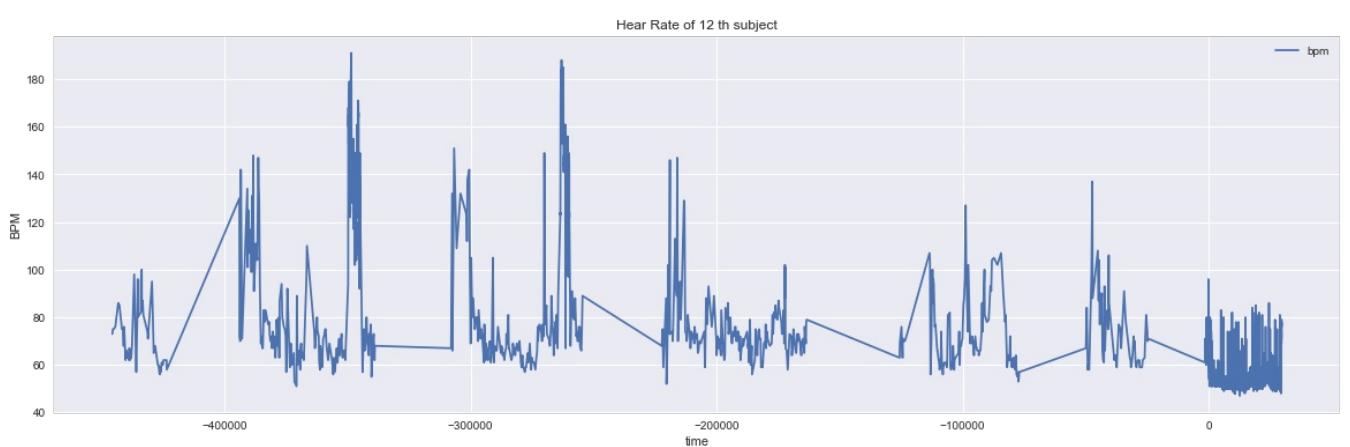
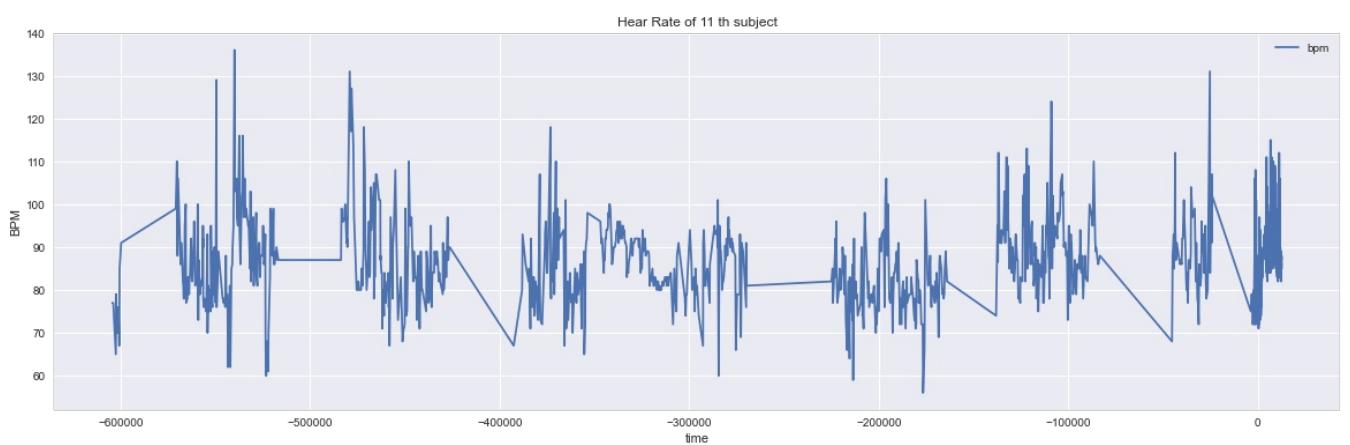
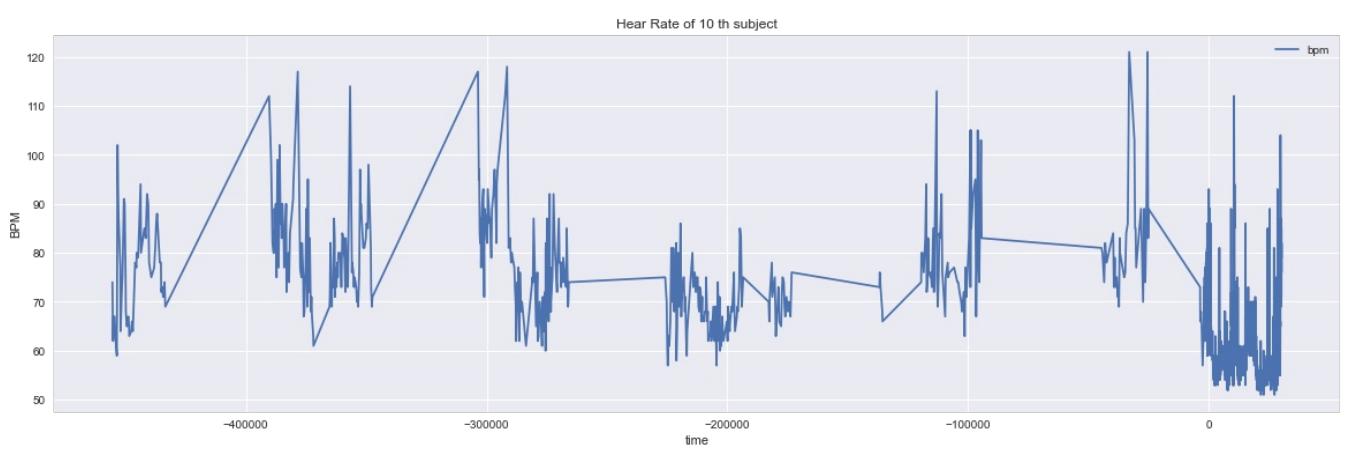
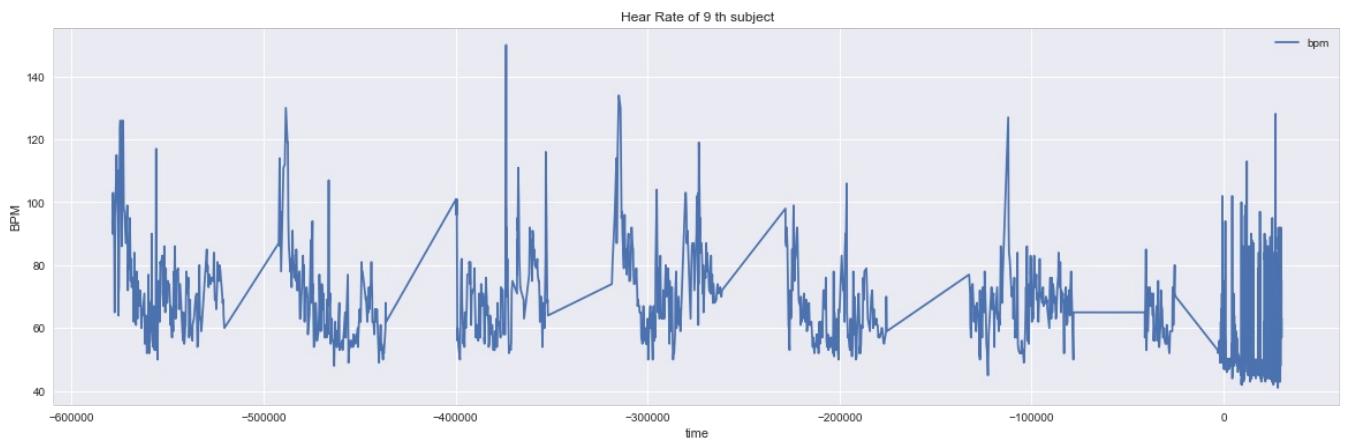
```
# 각 피험자별 heart rate 출력

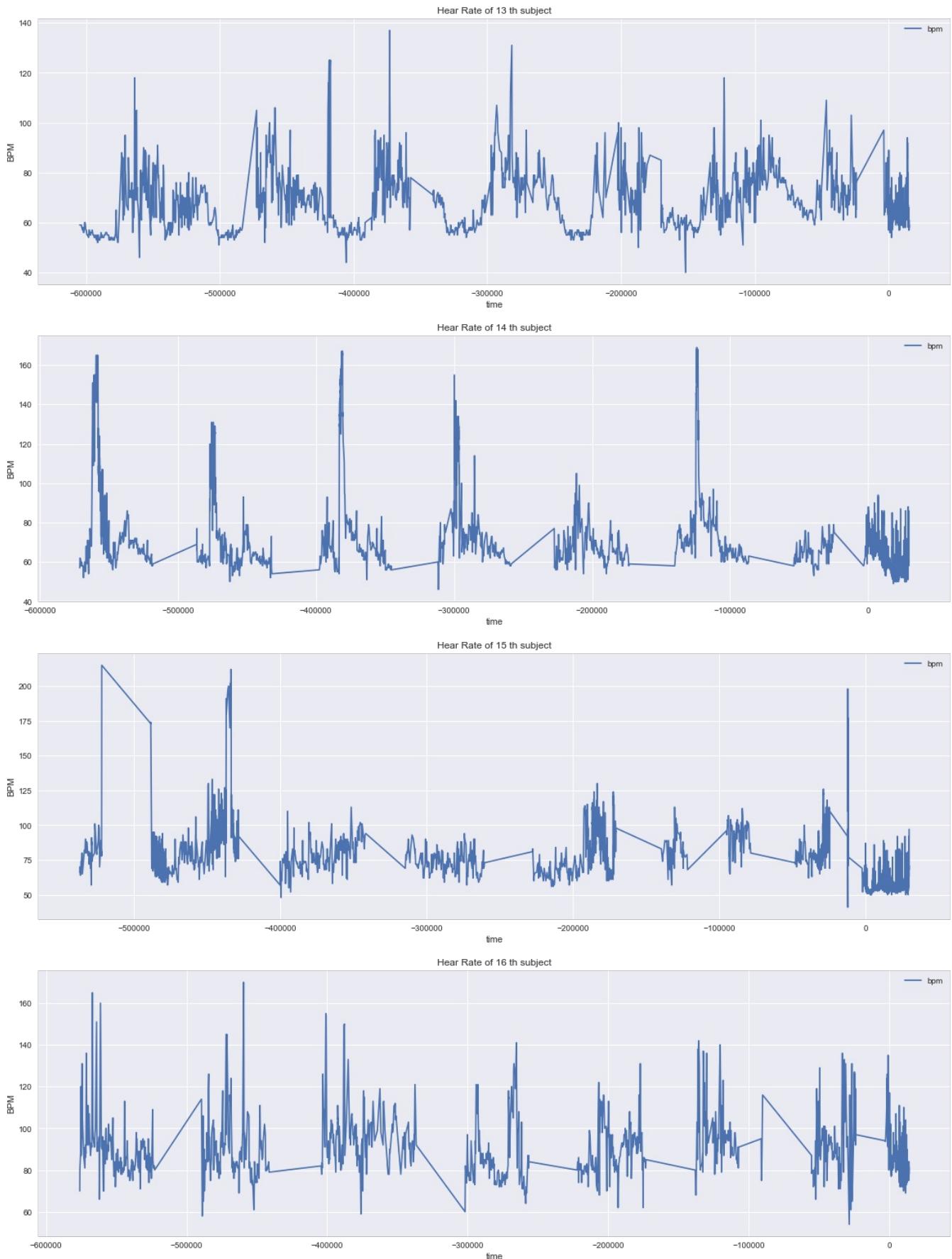
for i in range(len(heart_rate)):
    f, ax = plt.subplots(1, 1, figsize=(20, 6))
    plt.title(f'Hear Rate of {i} th subject')
    plt.xlabel('Time')
    plt.ylabel('BPM')
    heart_rate[i].plot(x='time', y='bpm', kind='line', ax=ax)
    plt.show()
```

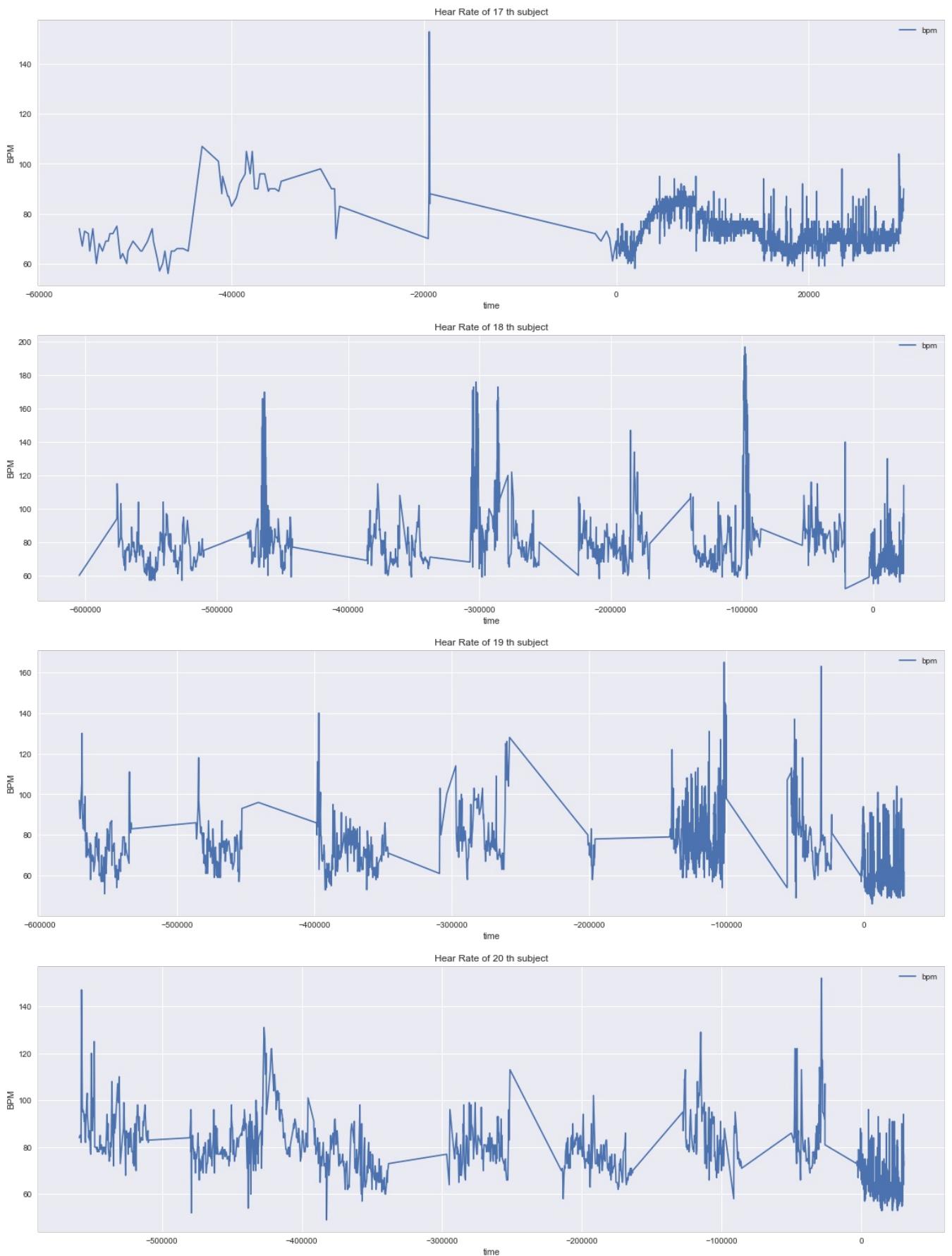


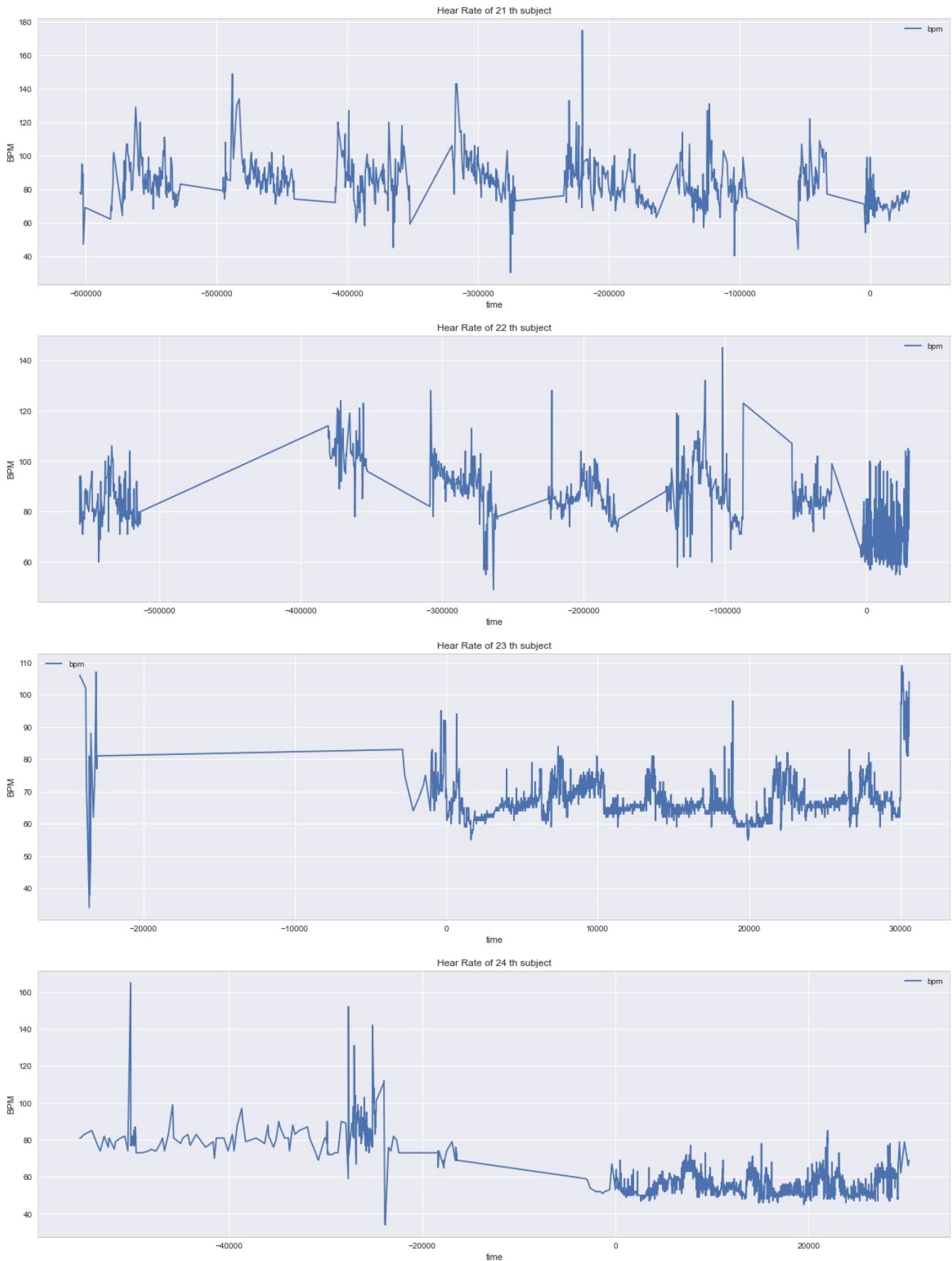


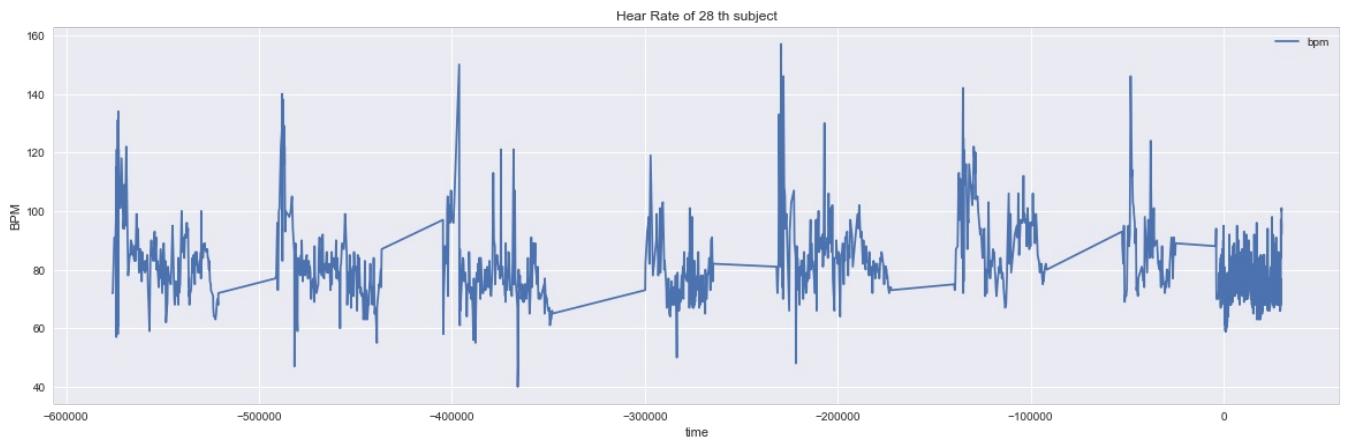
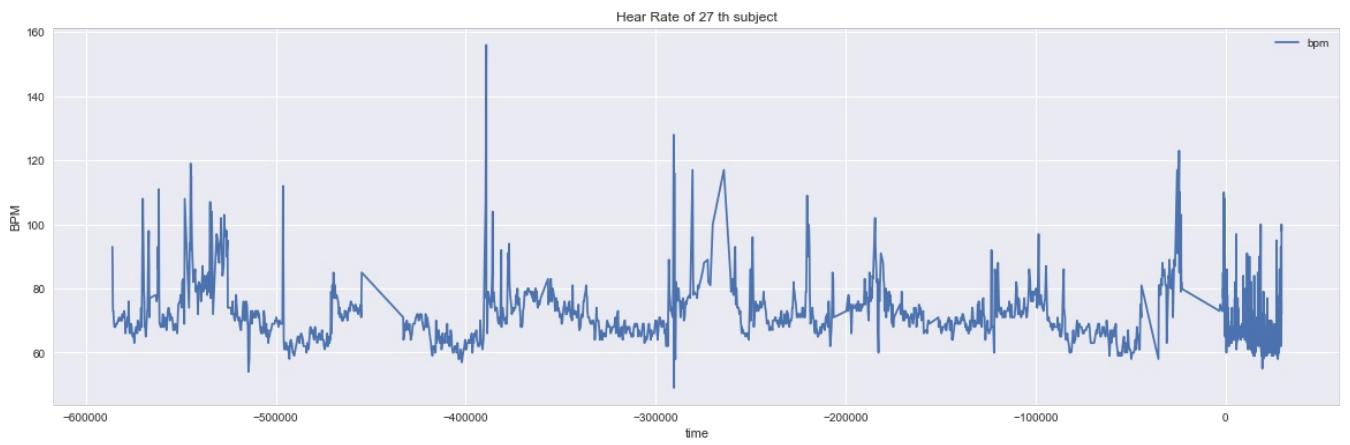
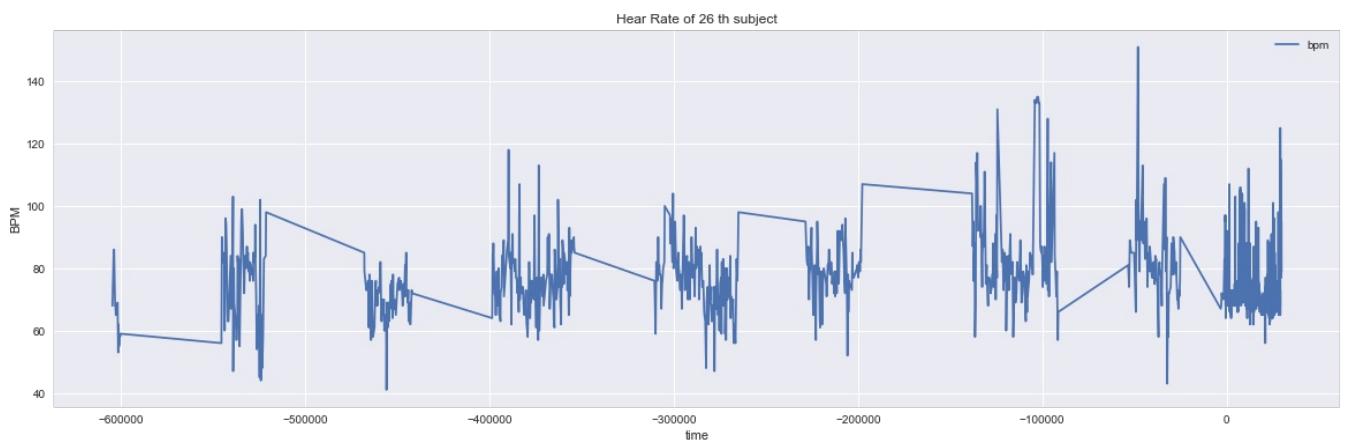
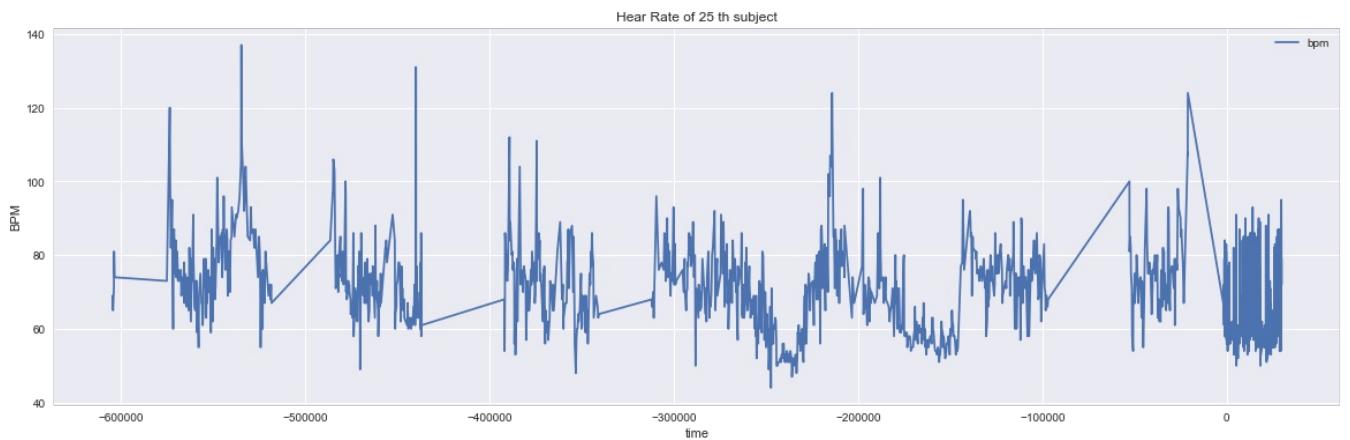


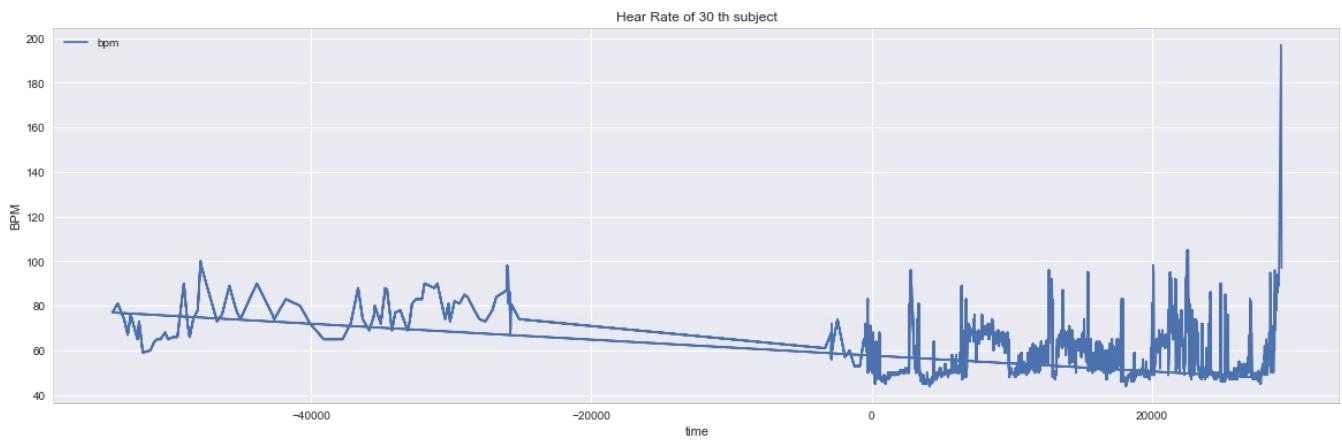
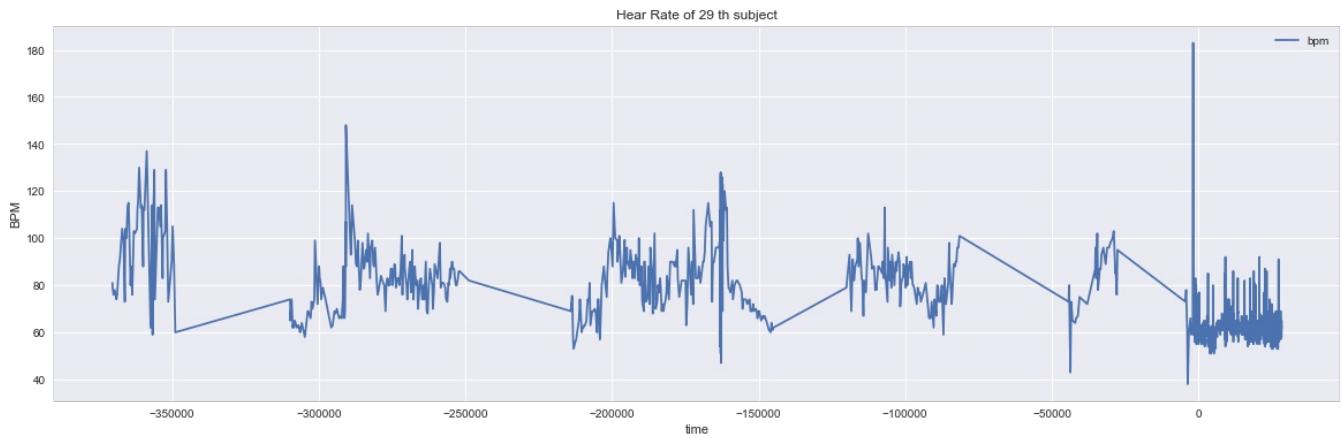






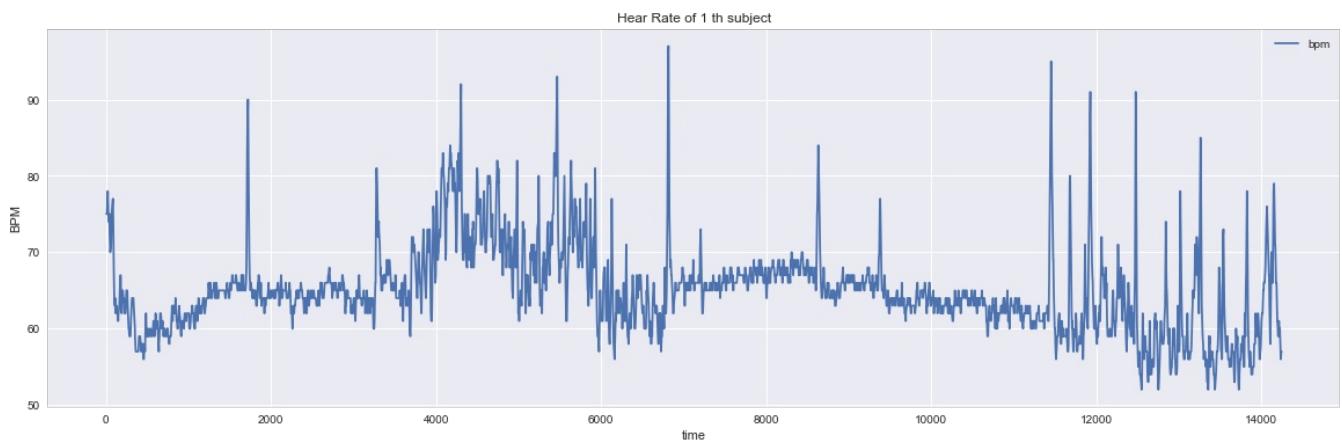
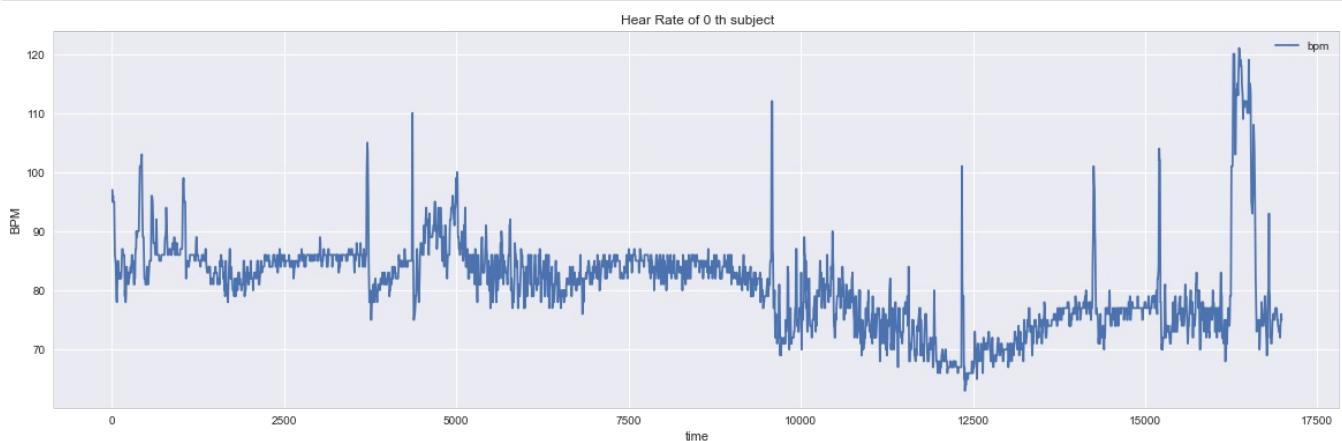


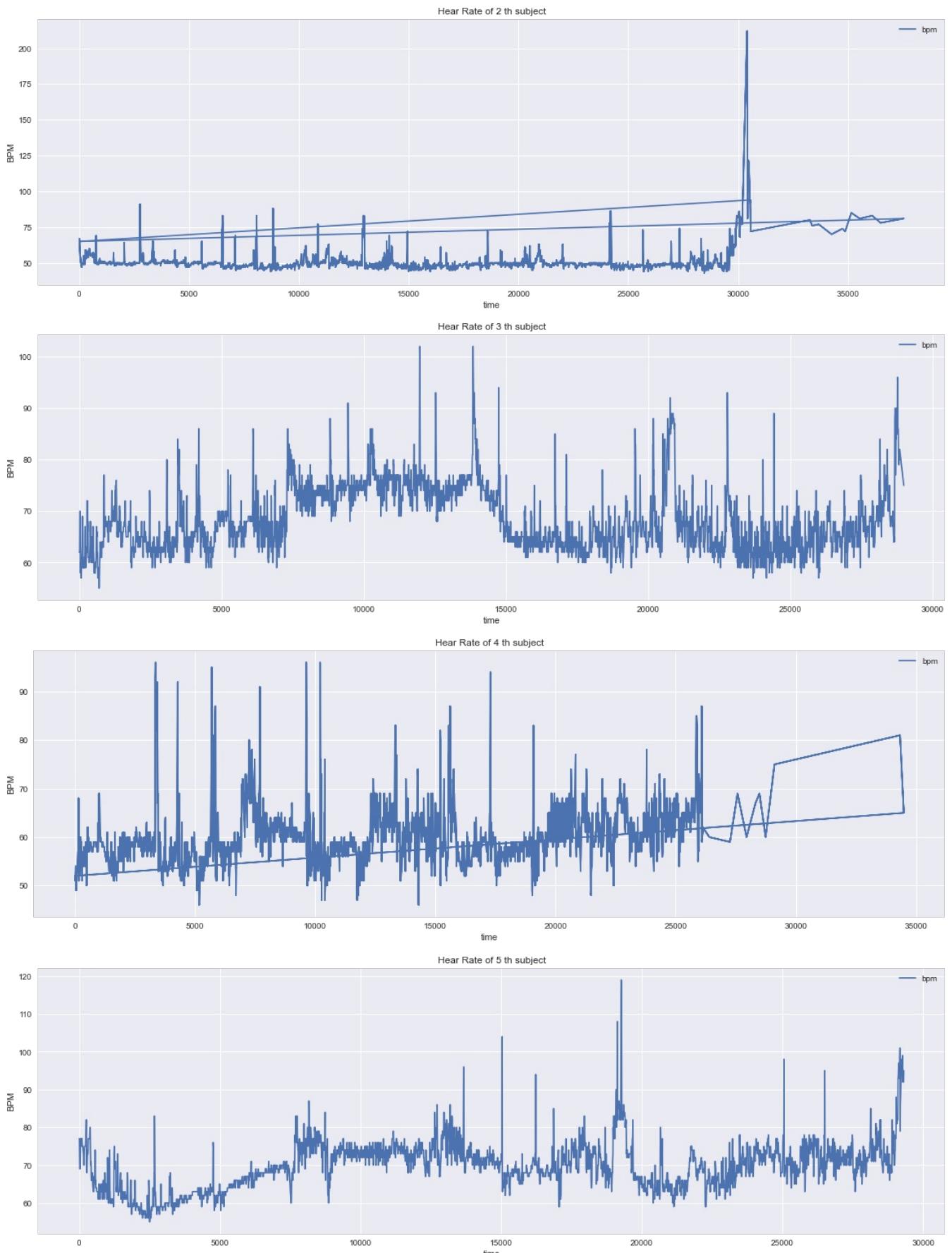


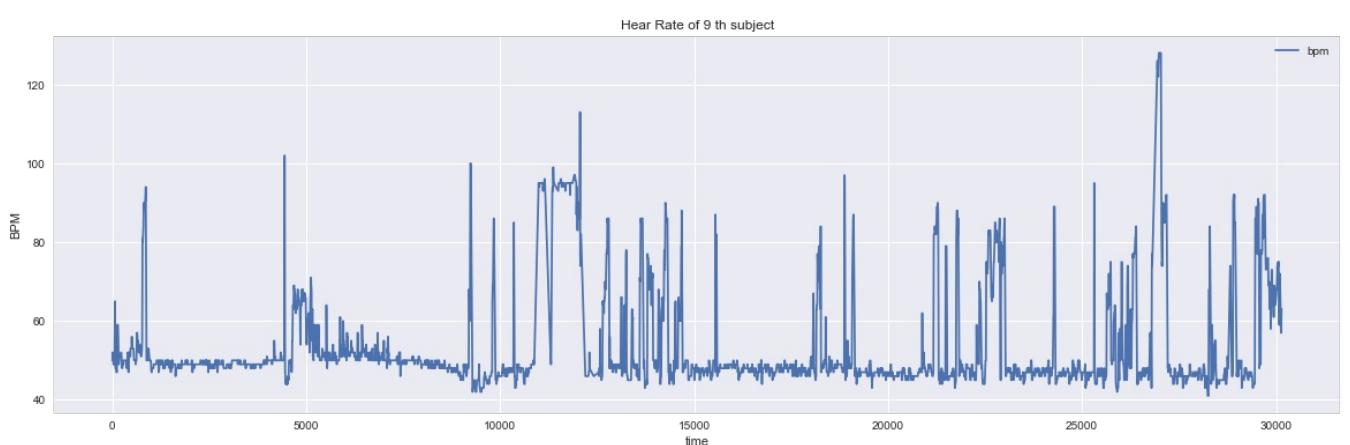
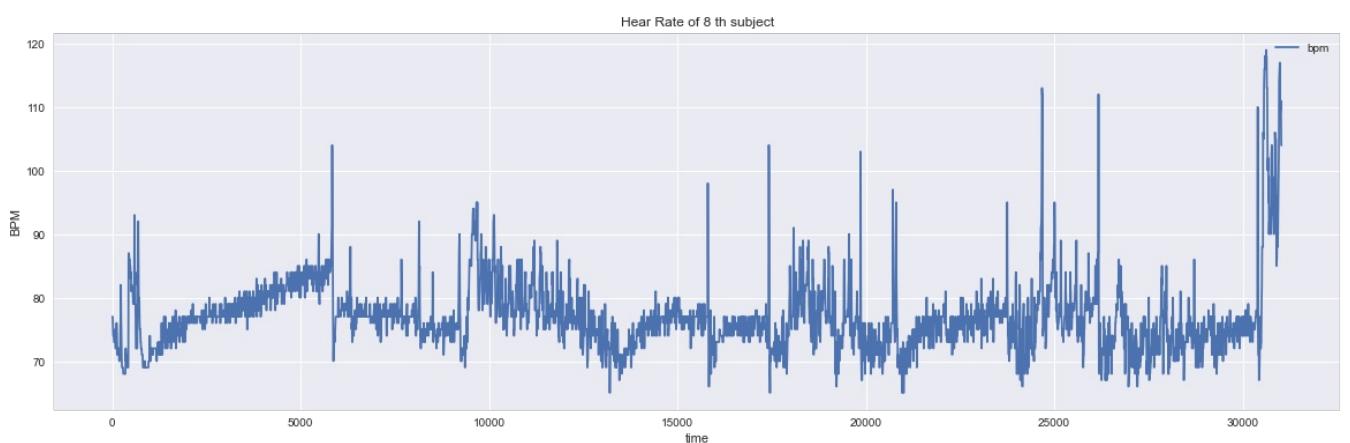
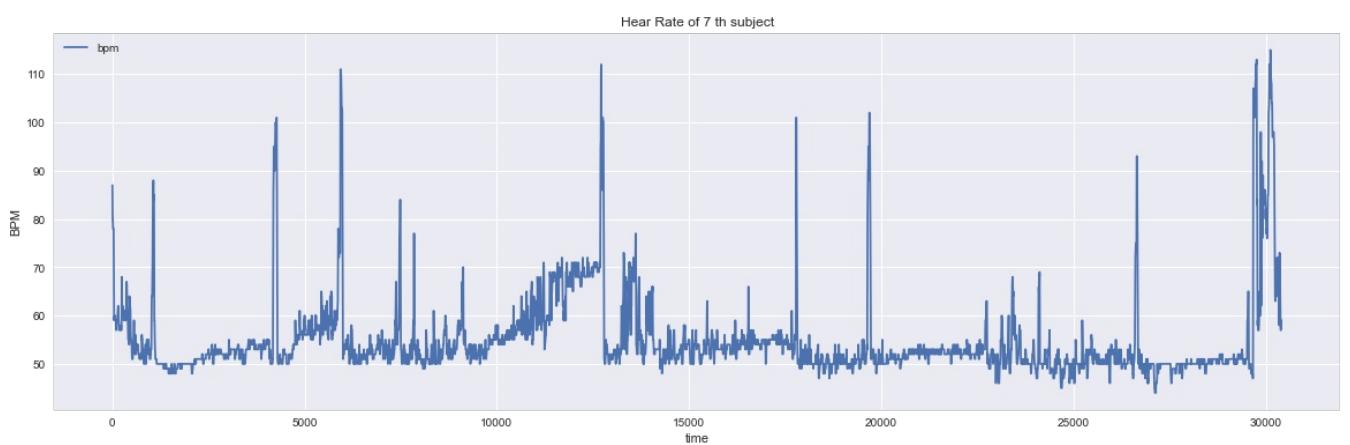
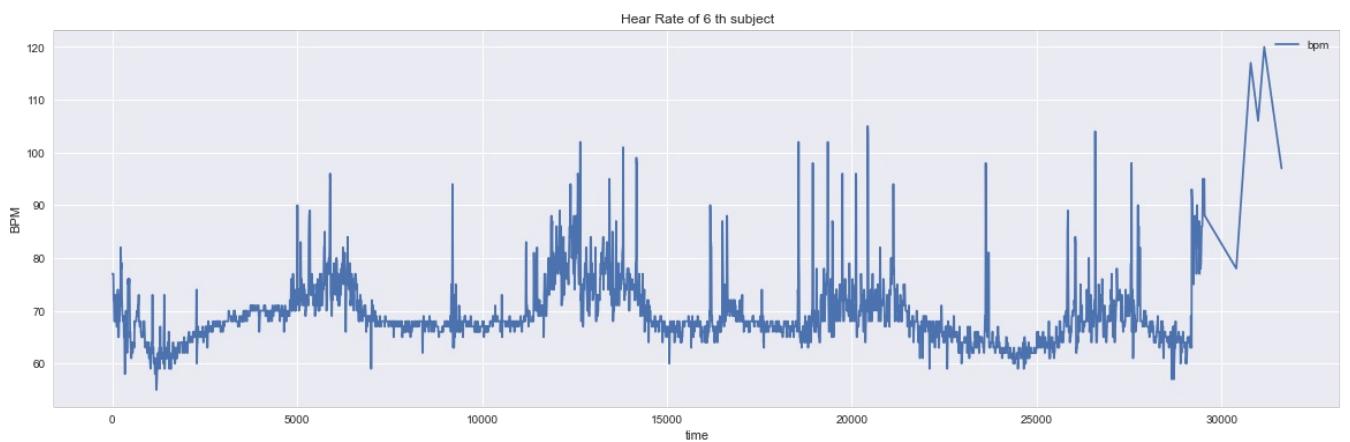


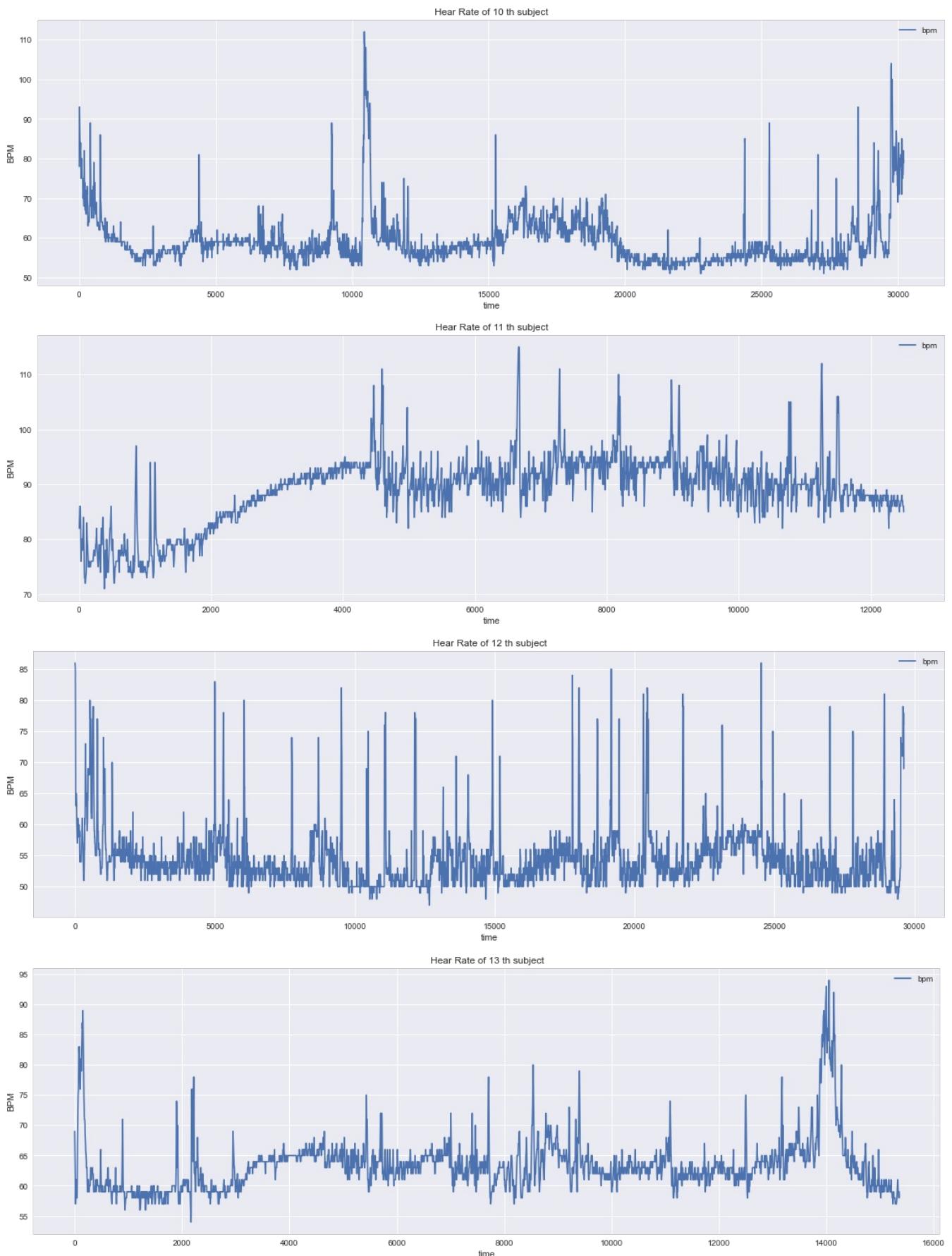
```
In [5]: # PSG 측정 시작 시간 이후의 heart rate 확인
```

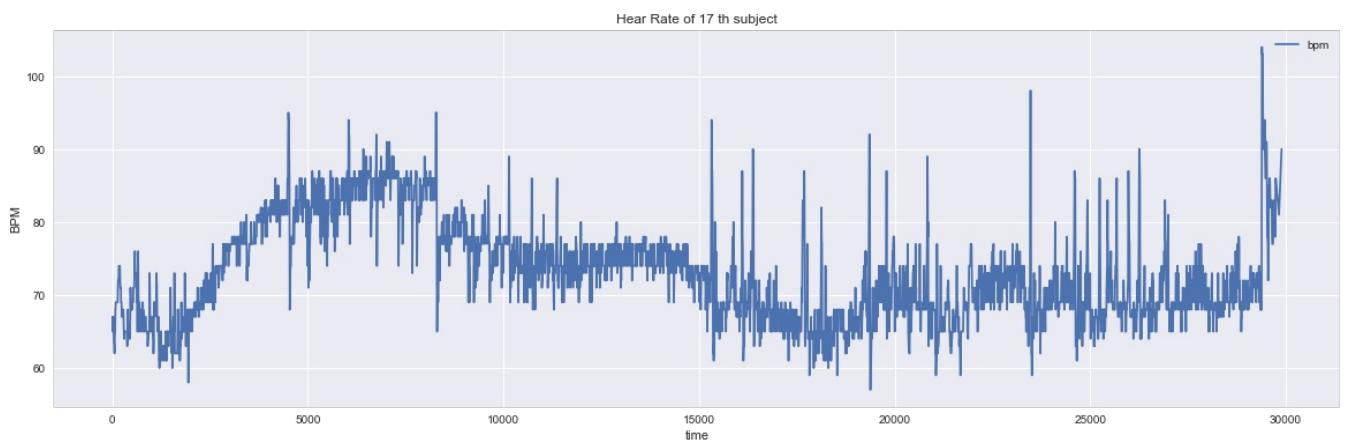
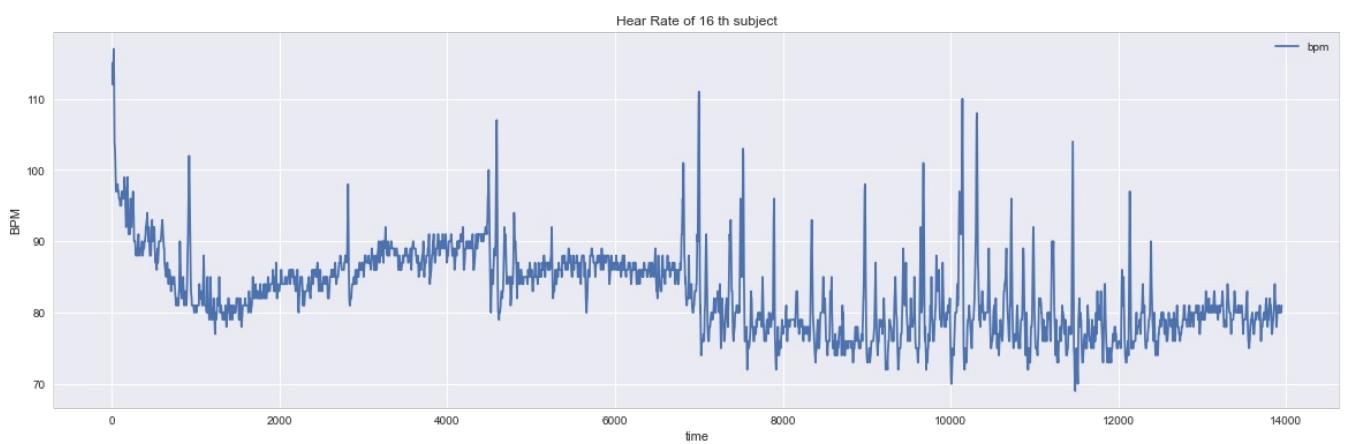
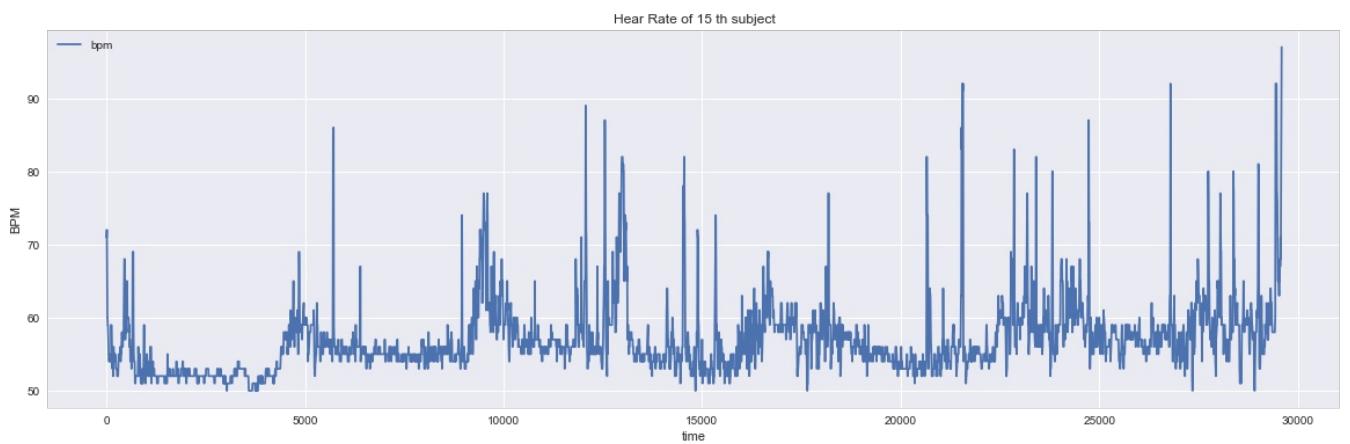
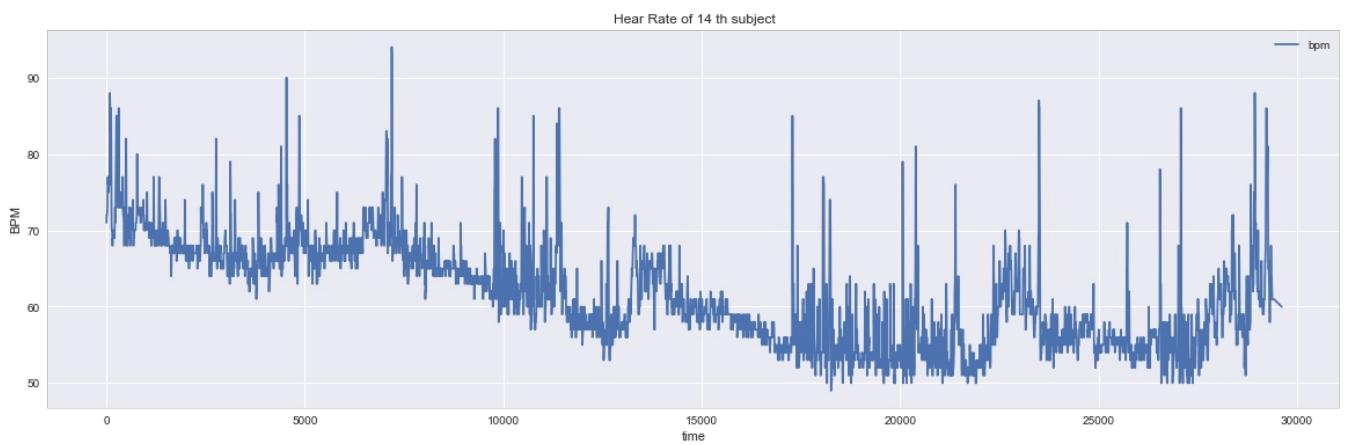
```
for i in range(len(heart_rate)):
    f, ax = plt.subplots(1, 1, figsize=(20, 6))
    plt.title(f'Hear Rate of {i} th subject')
    plt.xlabel('Time')
    plt.ylabel('BPM')
    heart_rate[i][heart_rate[i]['time'] >= 0].plot(x='time', y='bpm', kind='line', ax = ax)
    plt.show()
```

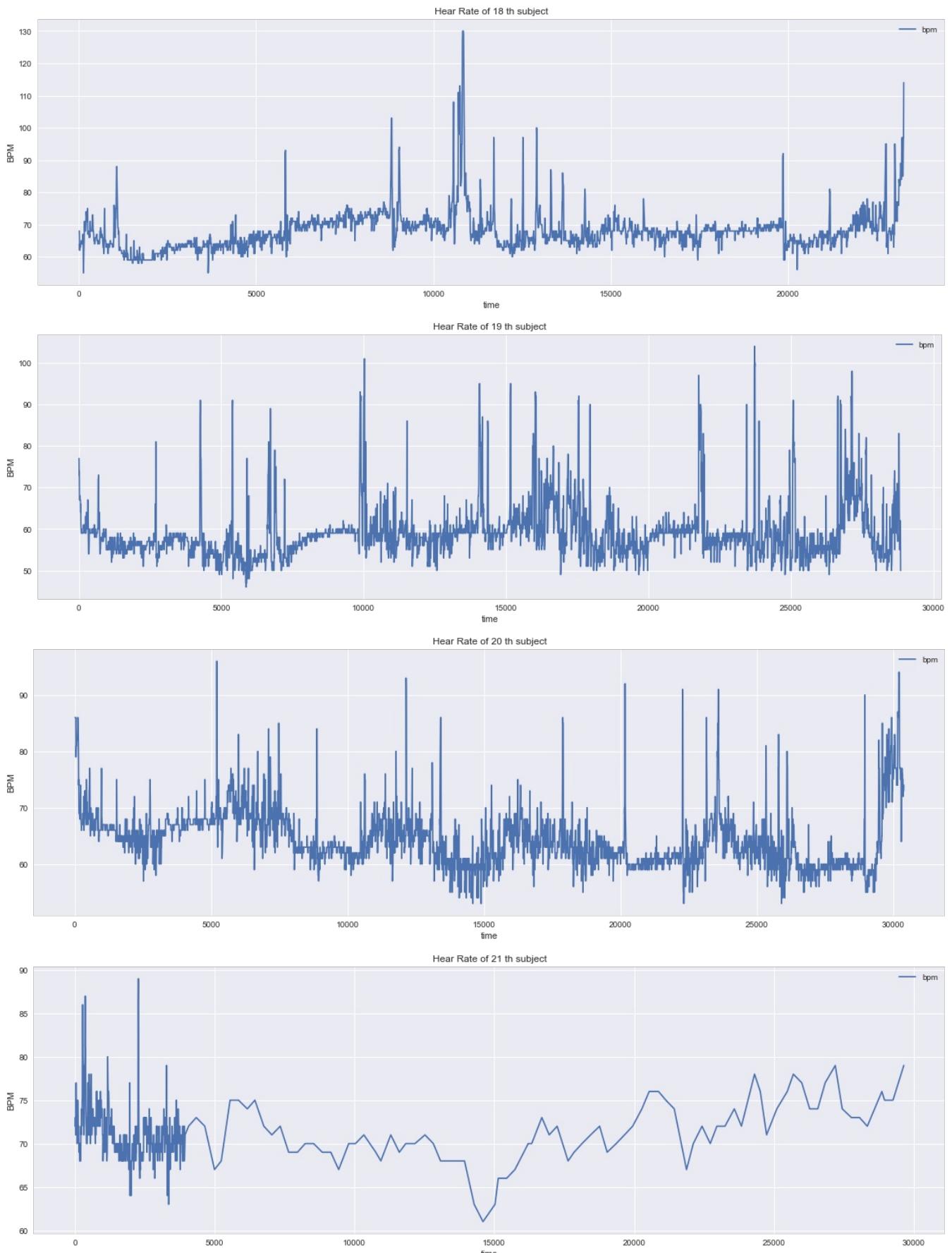


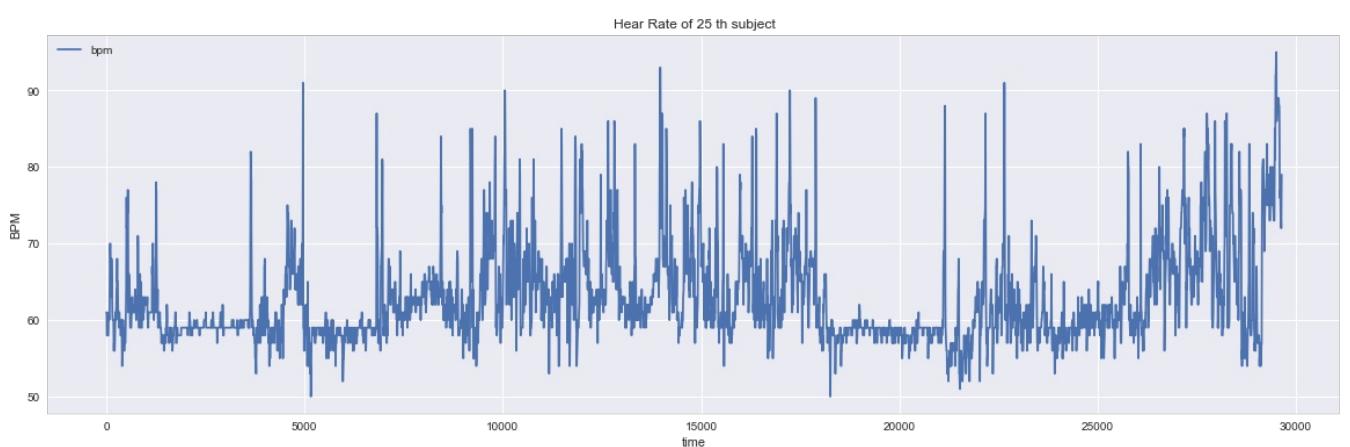
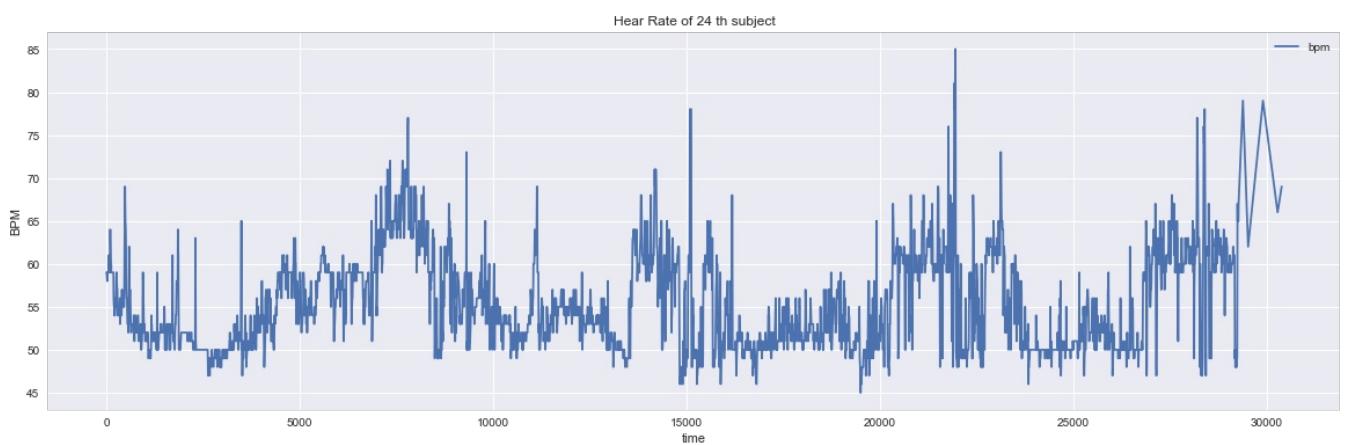
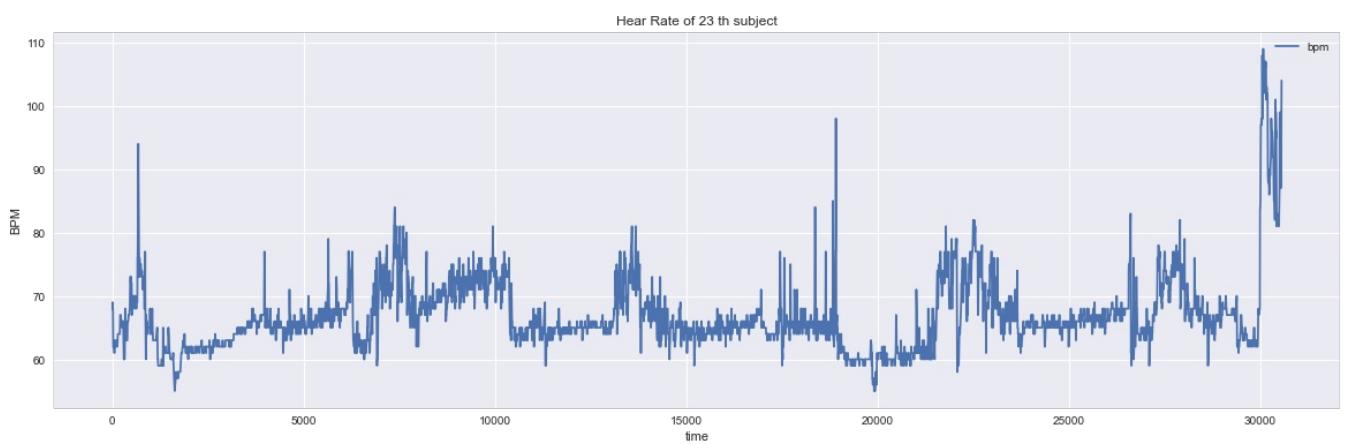
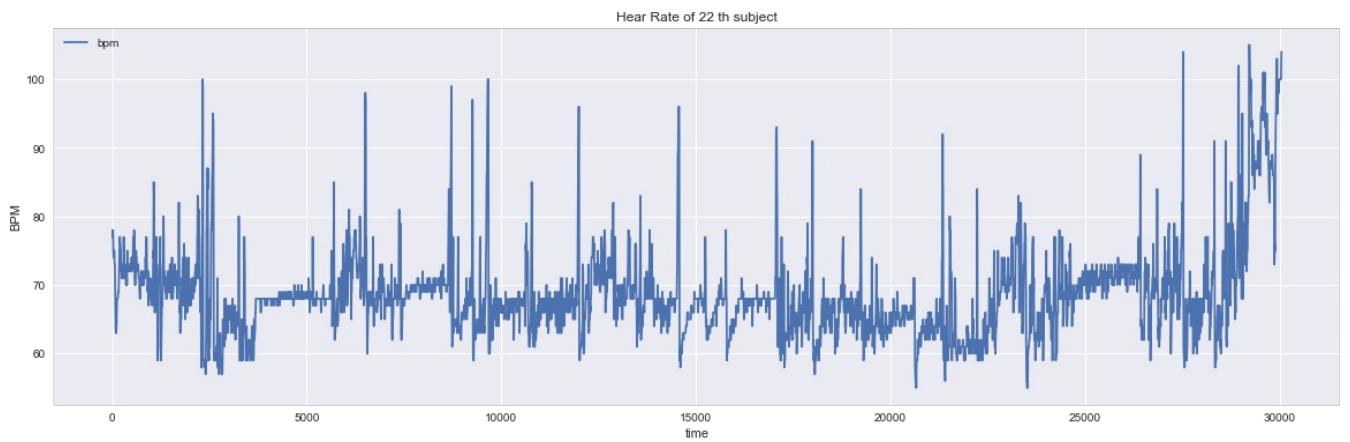


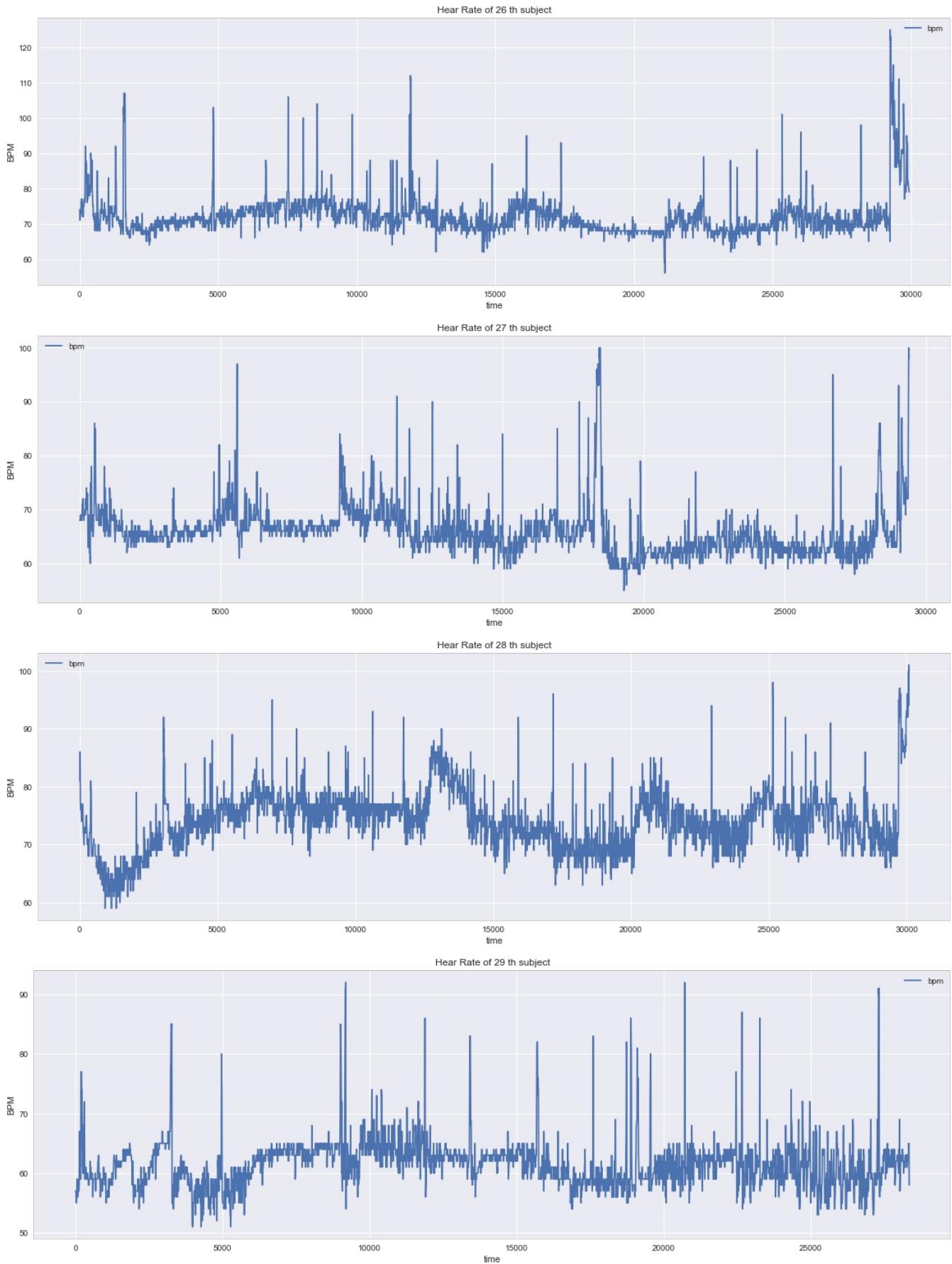


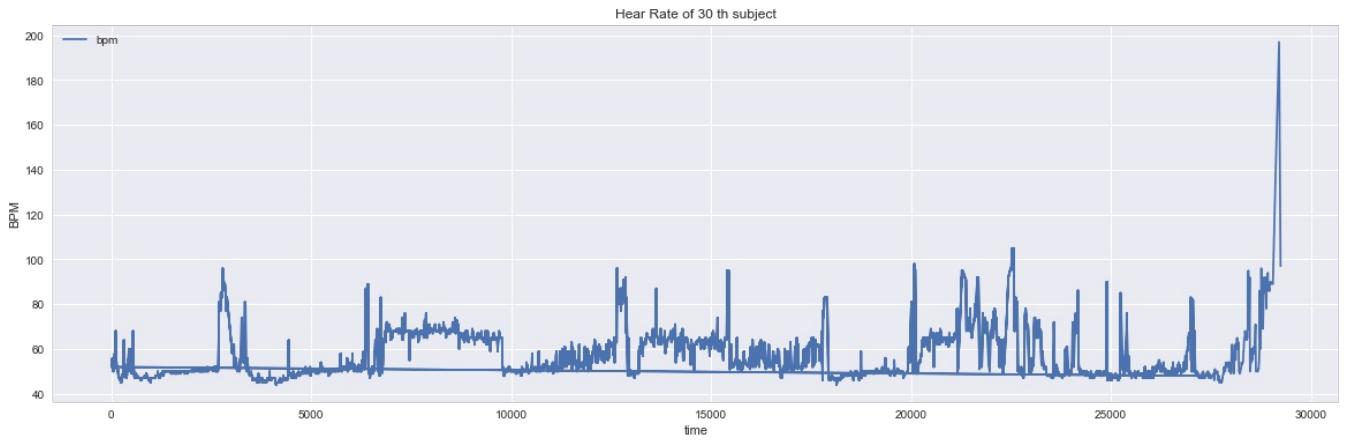












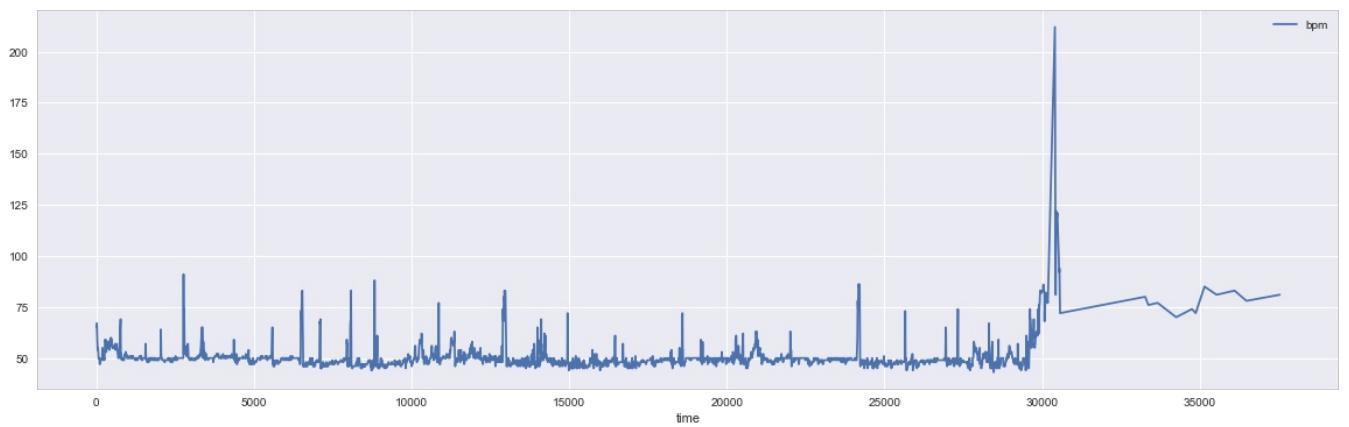
```
In [6]: # 2, 4, 30번 피험자는 값이 중복해서 나오므로 정리 필요 ( 두 번 나옴 )
df_list = [heart_rate[2], heart_rate[4], heart_rate[30]]
```

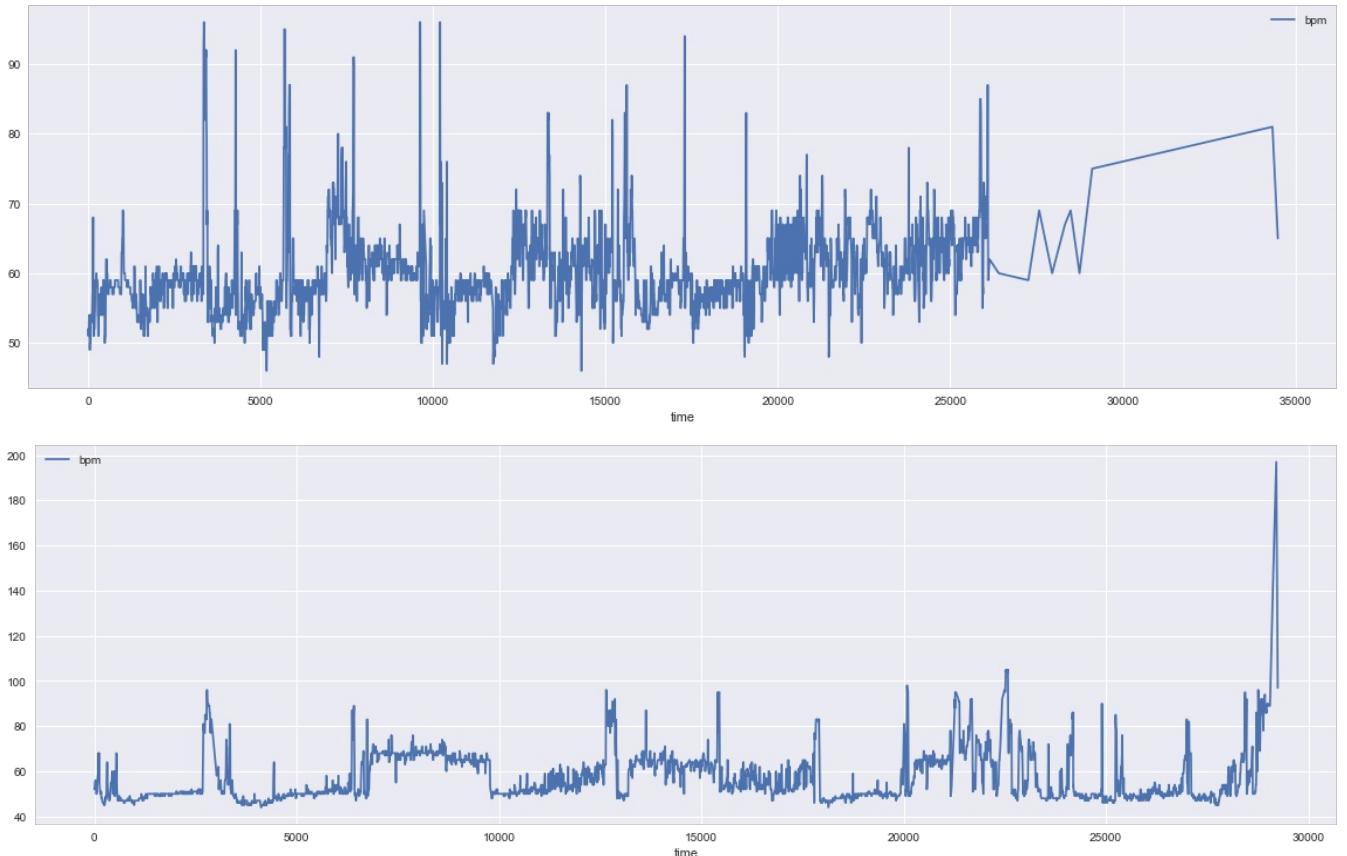
```
for df in df_list:
    # 데이터프레임의 인덱스와 Time 열을 순회하면서 감소하는 지점을 찾음
    prev_val = -999999
    time_idx = []

    for index, row in df.iterrows():
        if row['time'] < prev_val:
            prev_val = row['time']
            time_idx.append(index)

    prev_val = row['time']

df2 = df.iloc[max(time_idx):, :]
f, ax = plt.subplots(1, 1, figsize=(20, 6))
df2[df2['time'] >=0].plot(x='time', y='bpm', kind='line', ax = ax)
plt.show()
```





Steps EDA

수면 이전의 데이터이므로 수면 분류를 할 때 사용하지 않을 것임

```
In [9]: # steps 데이터 프레임 리스트
steps = []

dir_name = ''
path = ''
file_list = os.listdir(path)

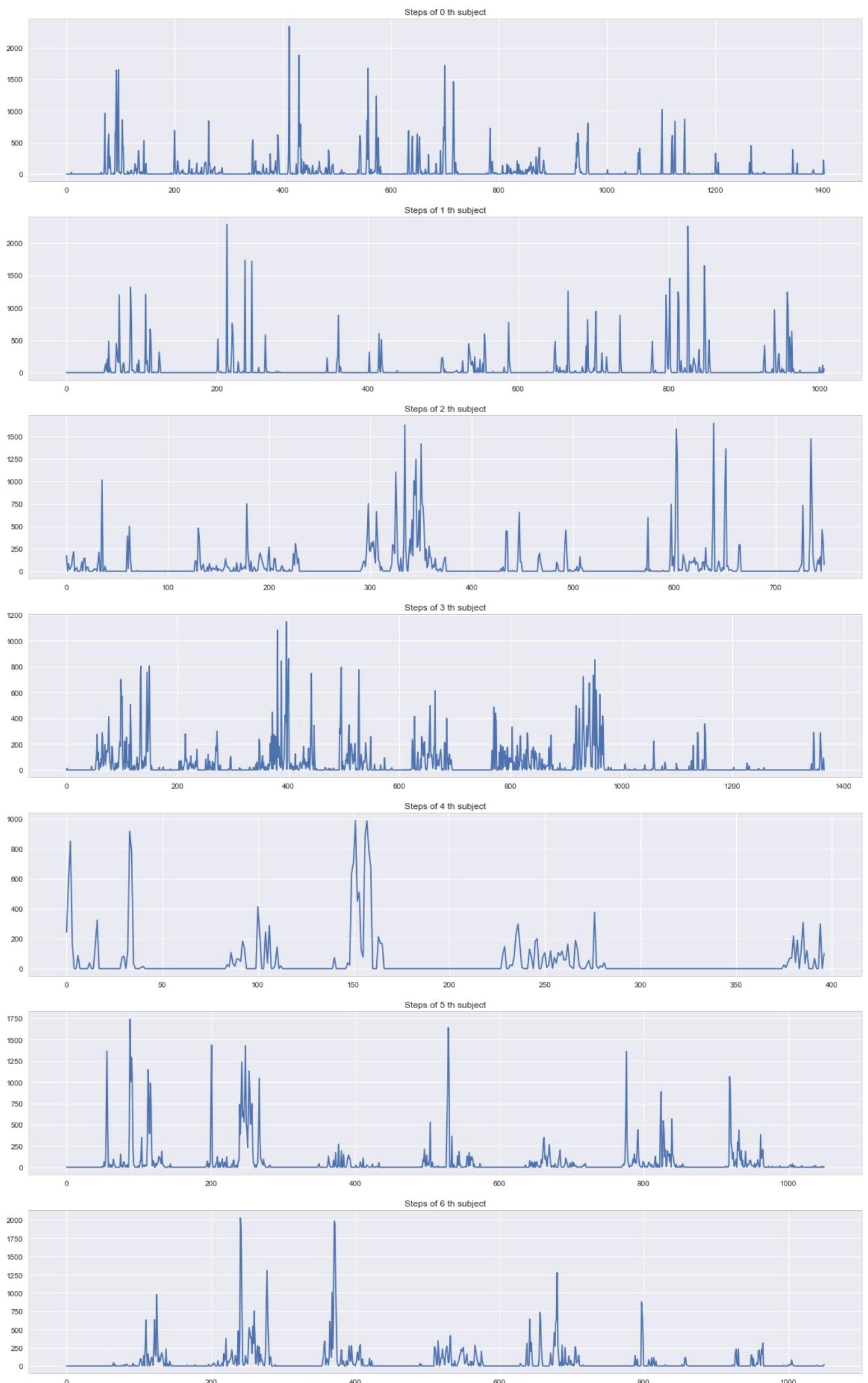
# 주어진 파일 이름에서 숫자 부분을 추출
# 파일 이름에서 정수 값을 추출하여 파일을 숫자순으로 정렬하는 데 사용
def get_number(filename):
    return int(re.search(r'\d+', filename).group())

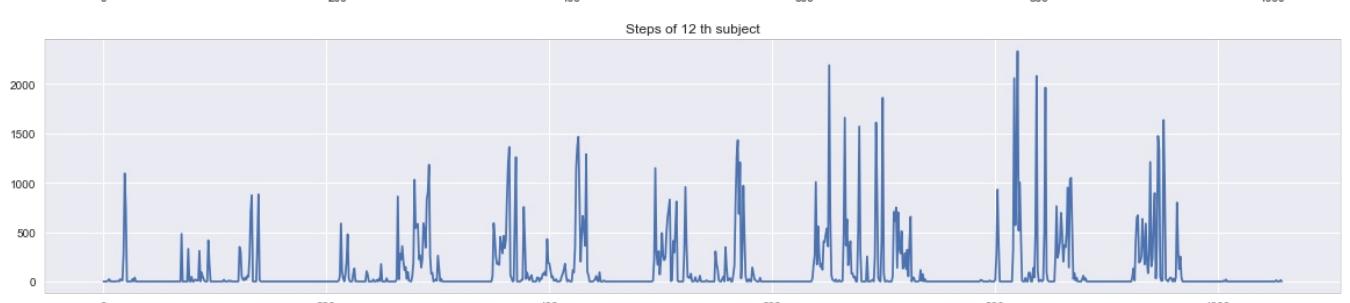
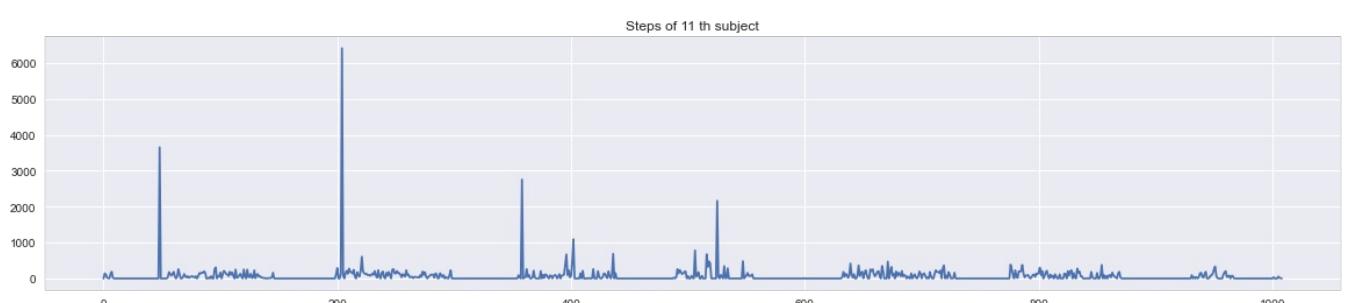
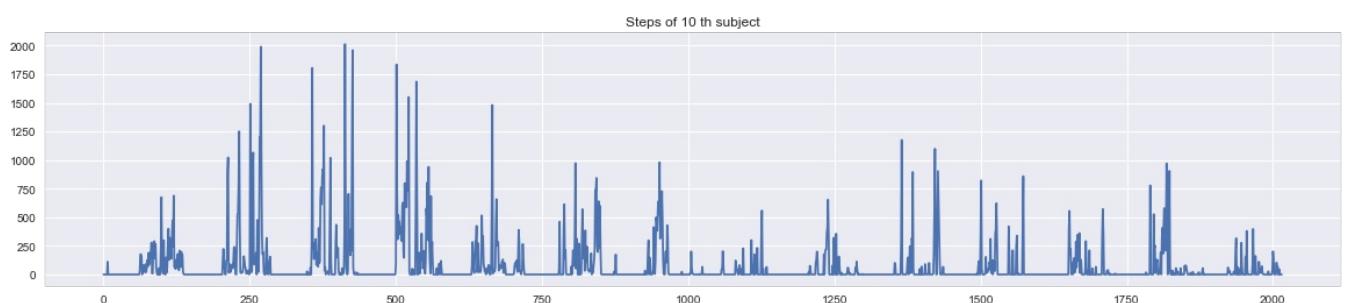
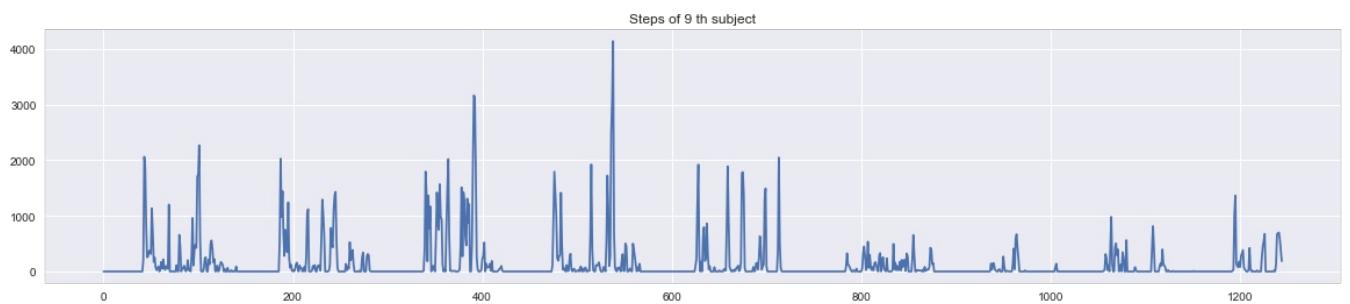
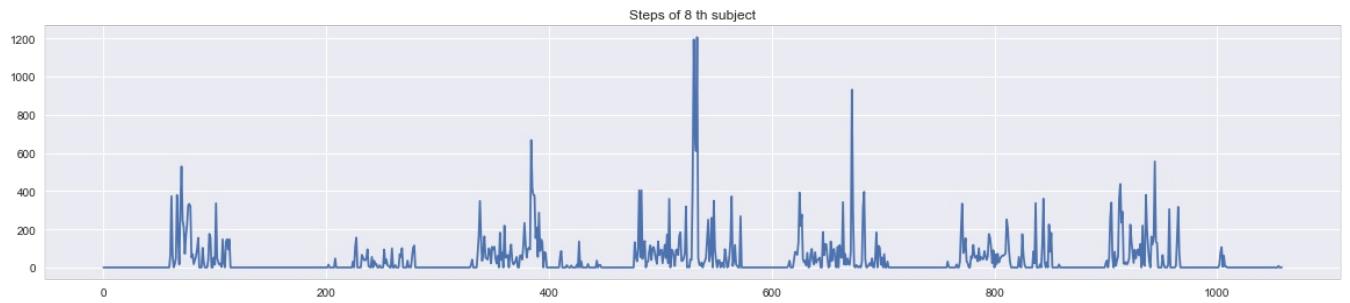
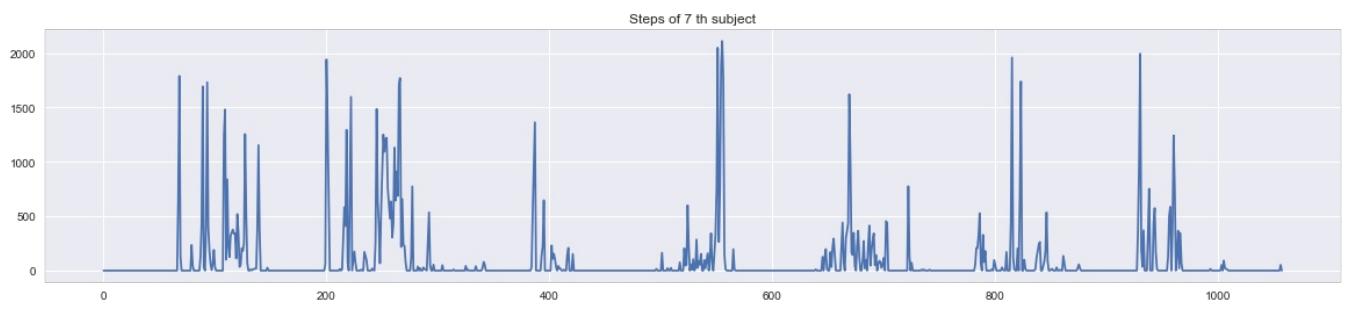
# 리스트 파일 이름 순으로 정렬
file_list_txt = sorted([file for file in file_list], key=get_number)

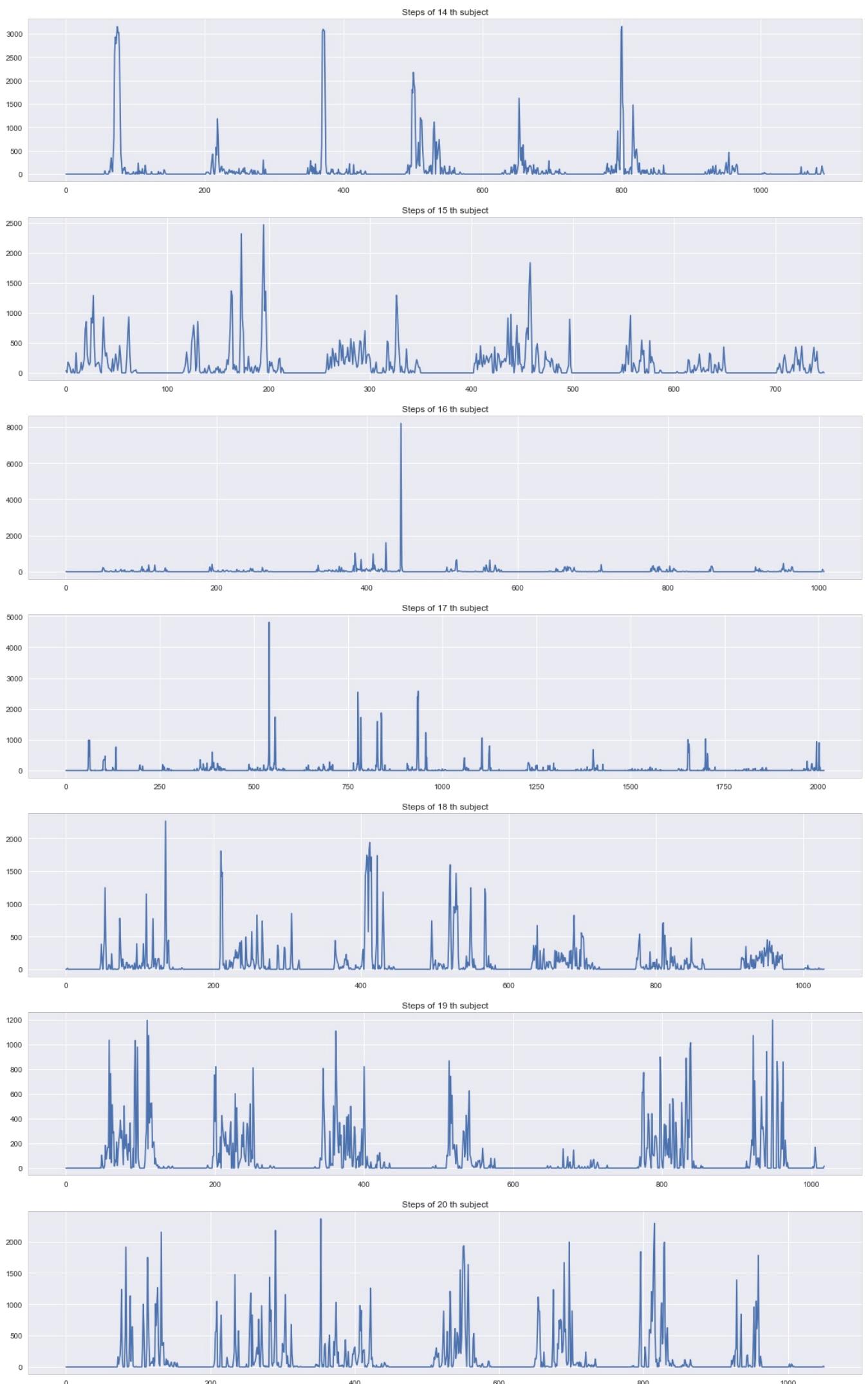
# 텍스트 파일을 csv 파일로 저장
for i in file_list_txt:
    df = pd.read_csv(dir_name+i, header=None, names=['time', 'step'])
    steps.append(df)
```

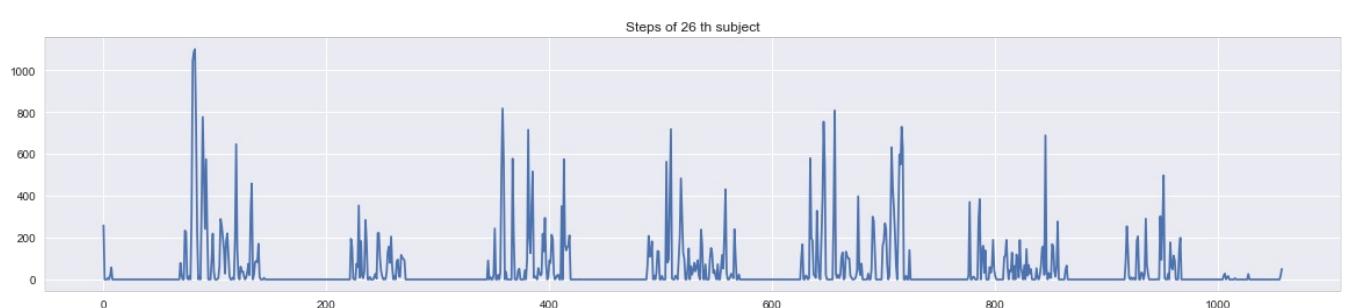
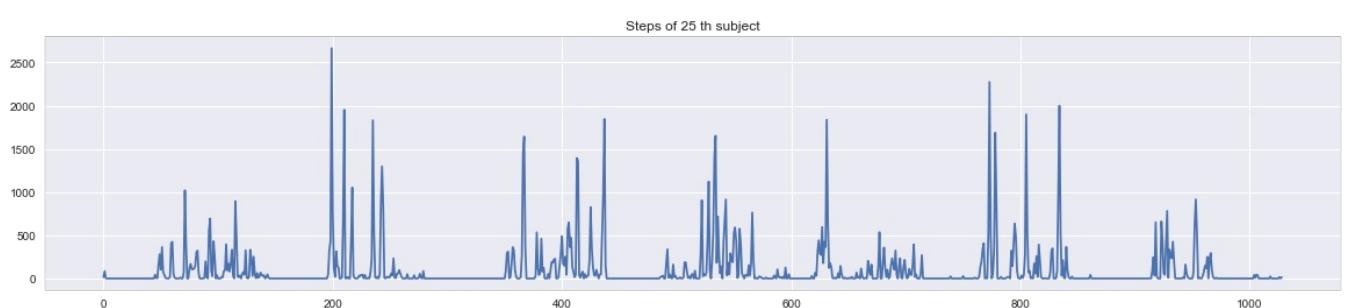
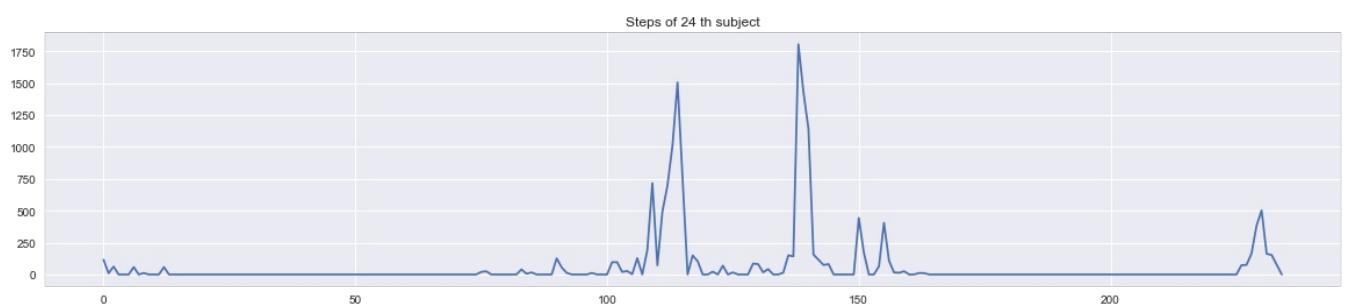
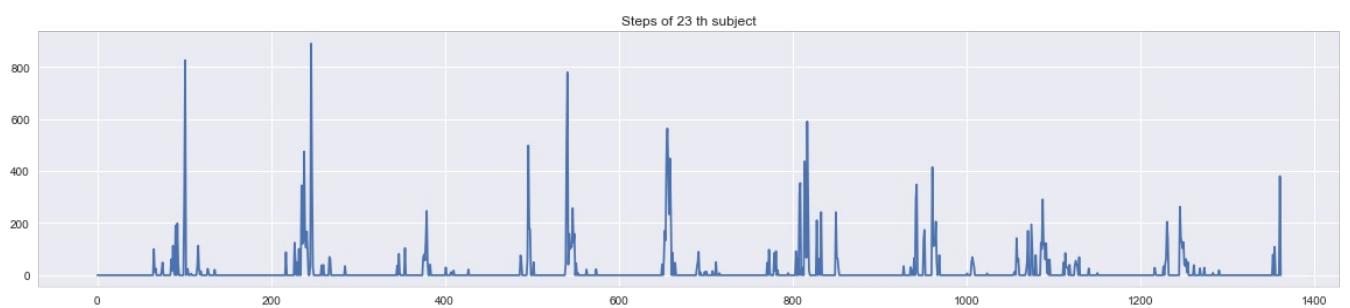
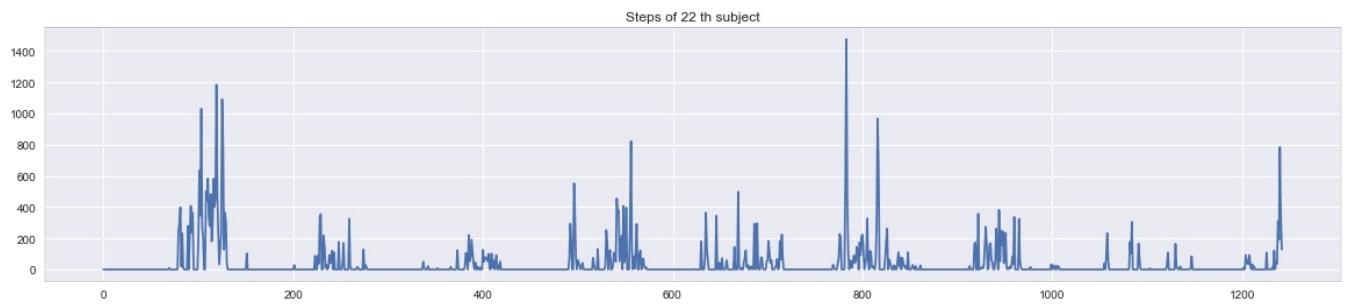
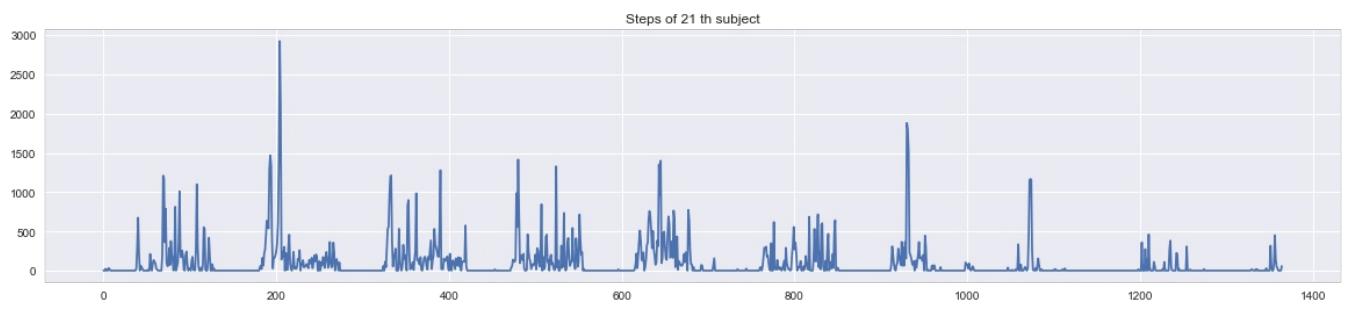
```
In [12]: # 각 피험자별 steps 출력

for i in range(len(steps)):
    plt.figure(figsize=(20, 4))
    plt.title(f'Steps of {i} th subject')
    steps[i]['step'].plot()
    plt.show()
```











Motion EDA

```
In [14]: # motion 데이터 프레임 리스트
motion = []

dir_name = ''
path = ''
file_list = os.listdir(path)

# 주어진 파일 이름에서 숫자 부분을 추출
# 파일 이름에서 정수 값을 추출하여 파일을 숫자순으로 정렬하는 데 사용
def get_number(filename):
    return int(re.search(r'\d+', filename).group())

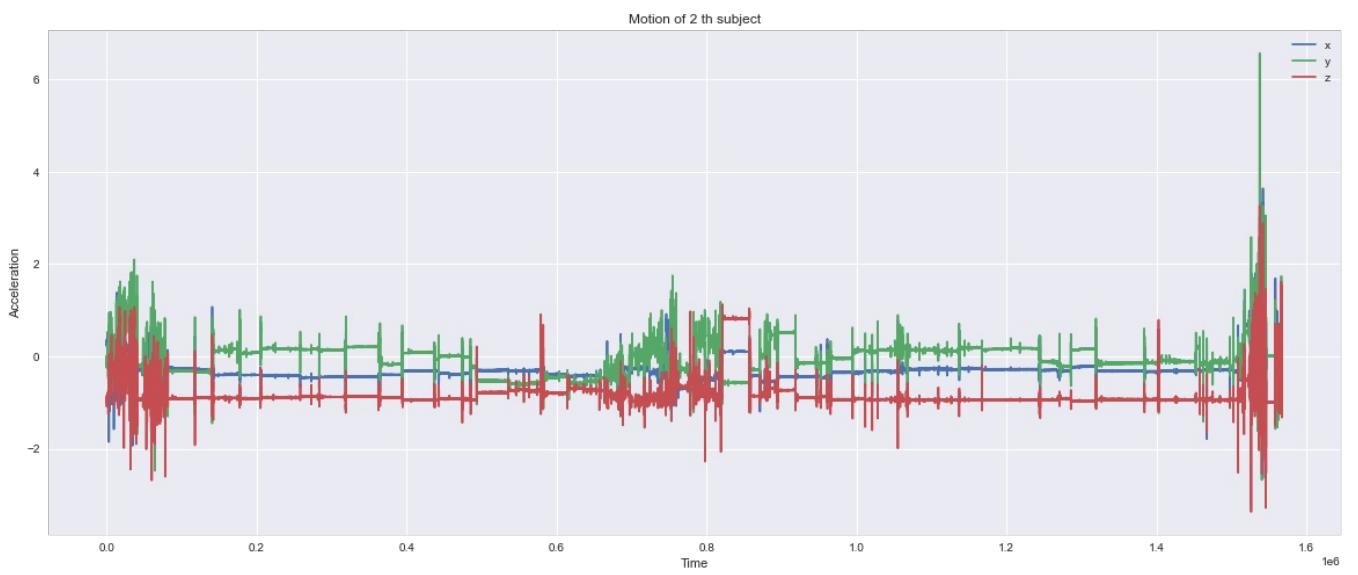
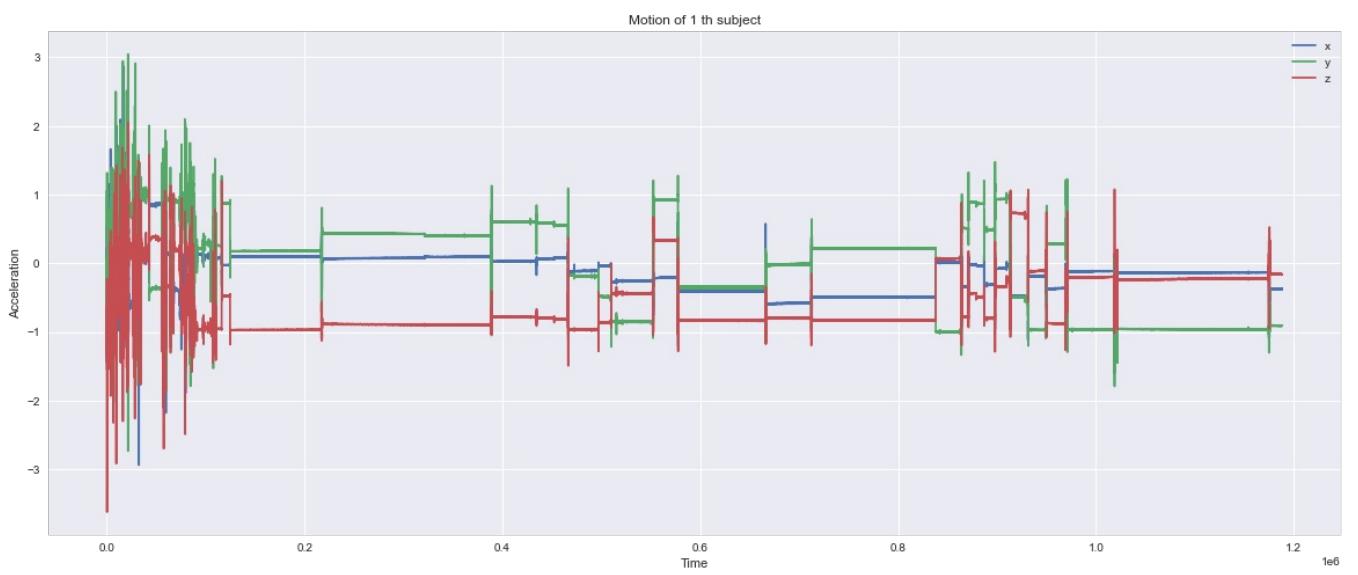
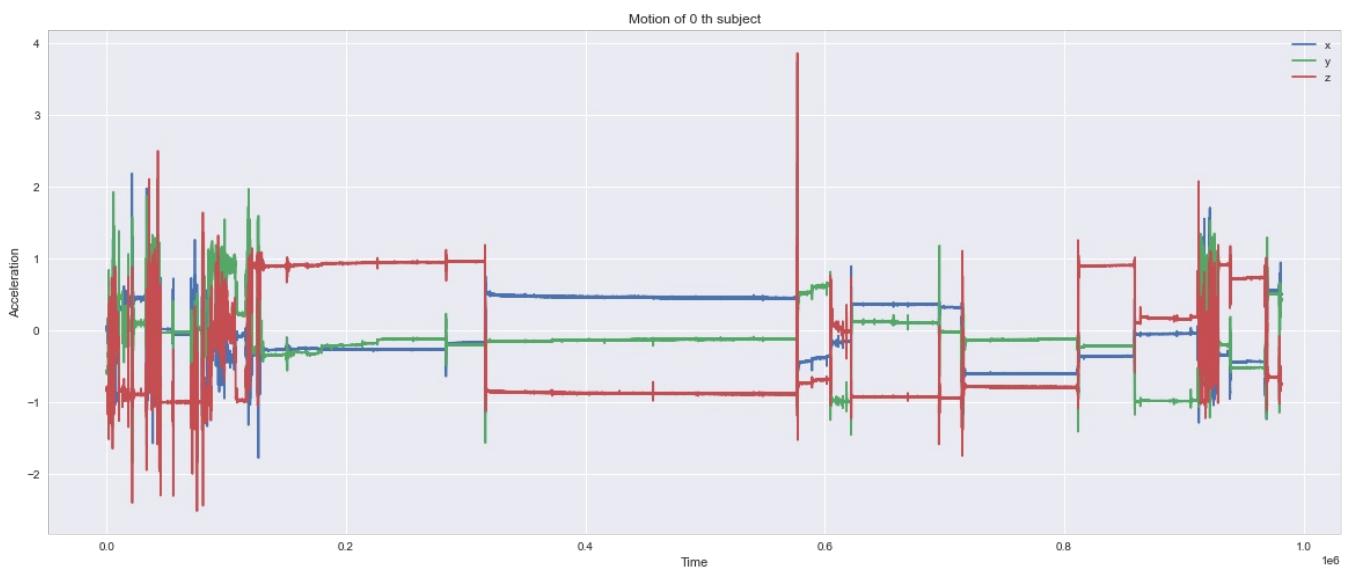
# 리스트 파일 이름 순으로 정렬
file_list_txt = sorted([file for file in file_list], key=get_number)

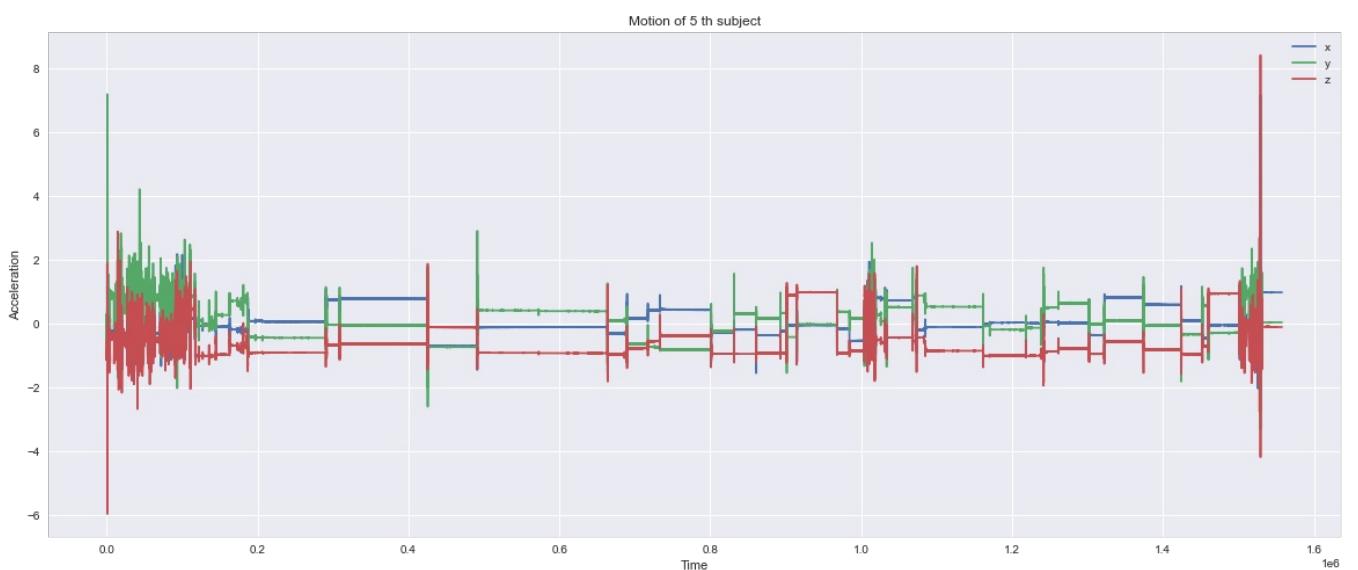
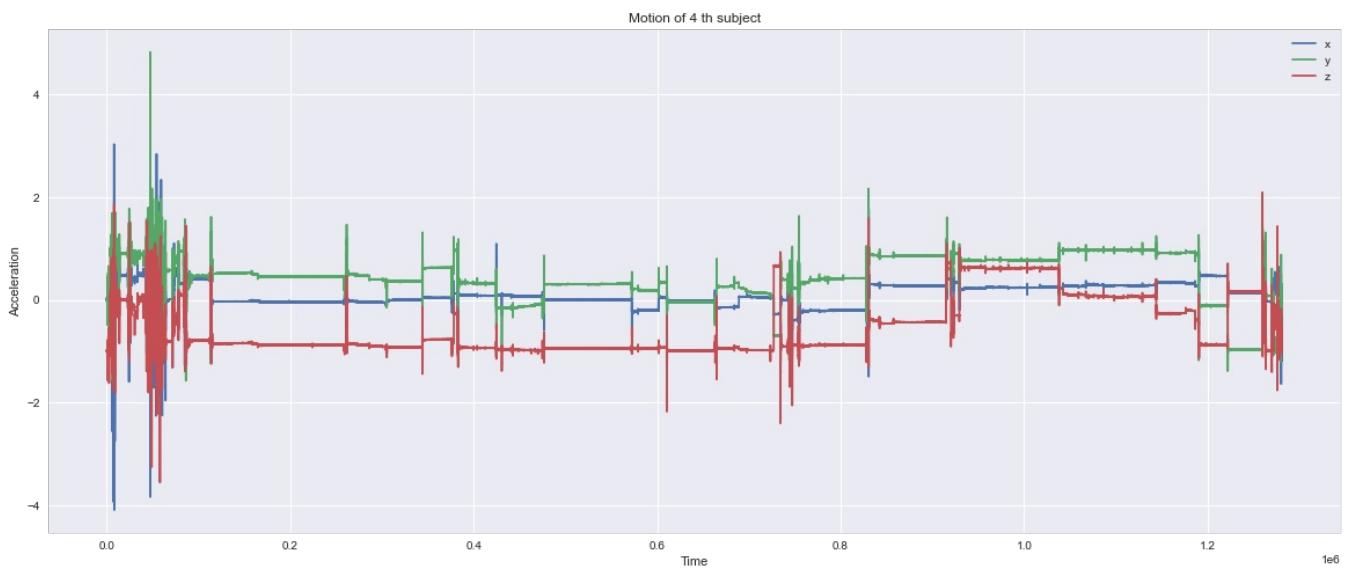
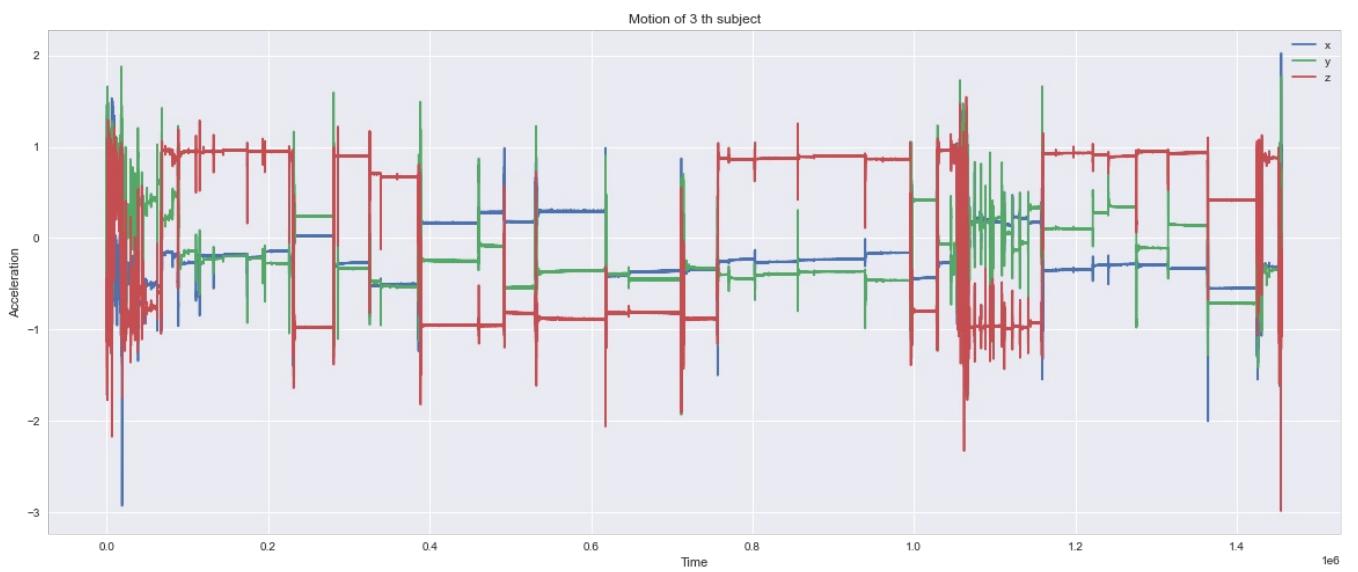
# 텍스트 파일을 csv 파일로 저장
for i in file_list_txt:
    df = pd.read_csv(dir_name+i, sep=' ', header=None, names=['time', 'x', 'y', 'z'])
    motion.append(df)
```

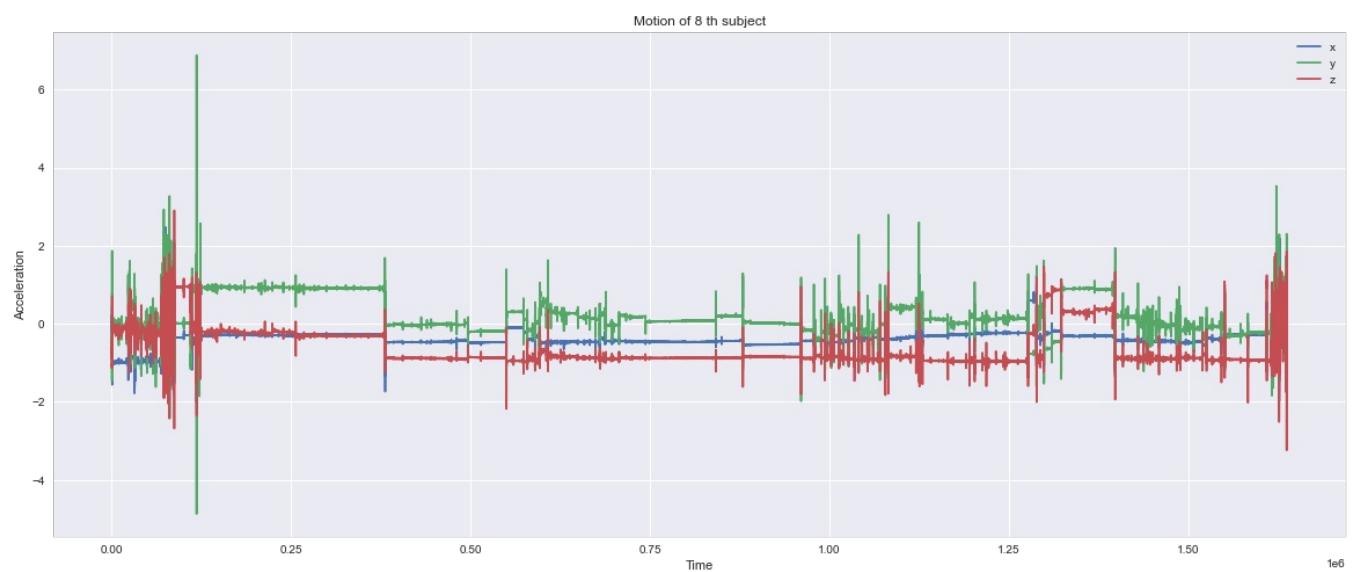
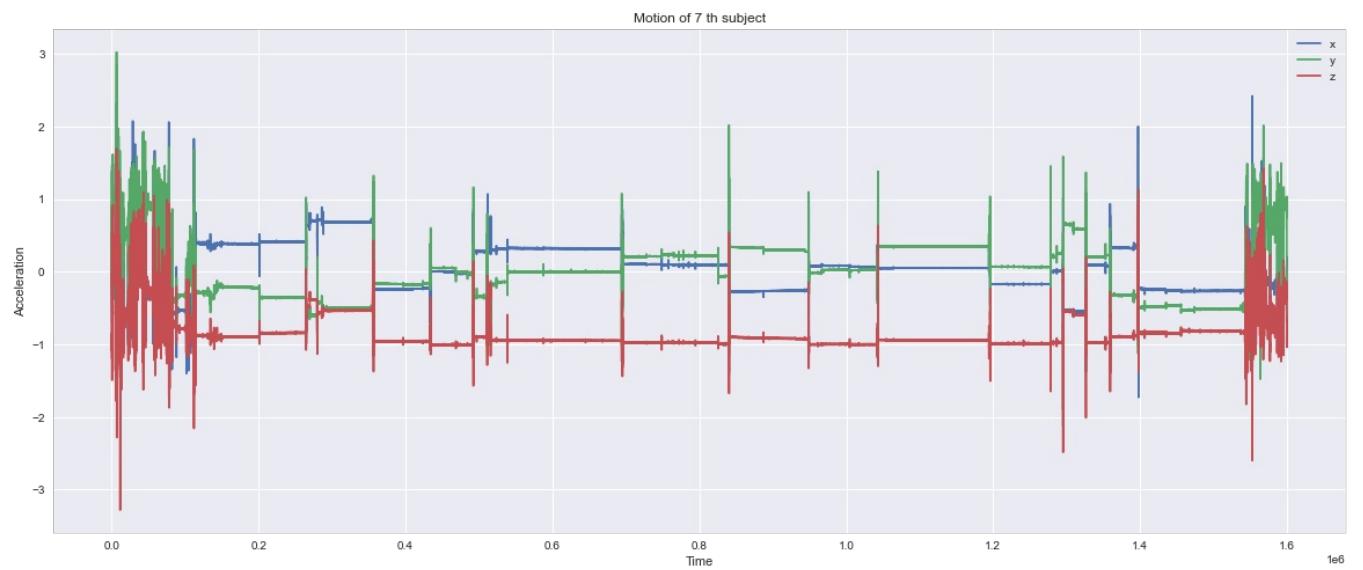
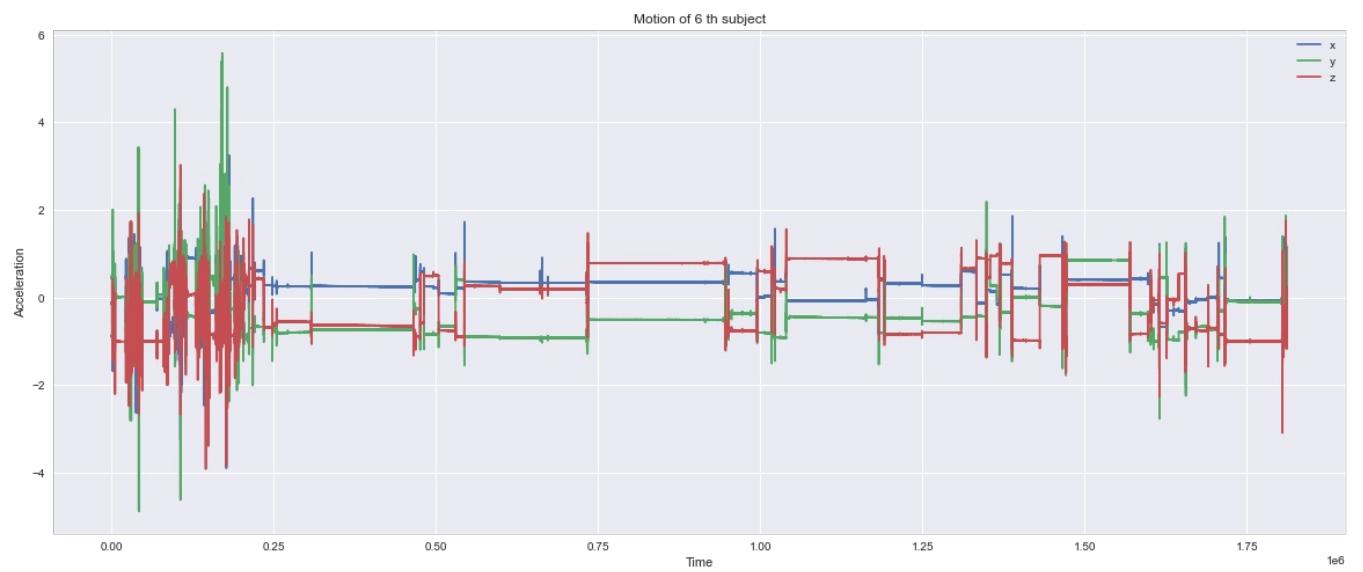
```
In [15]: # 각 피험자별 motion 출력

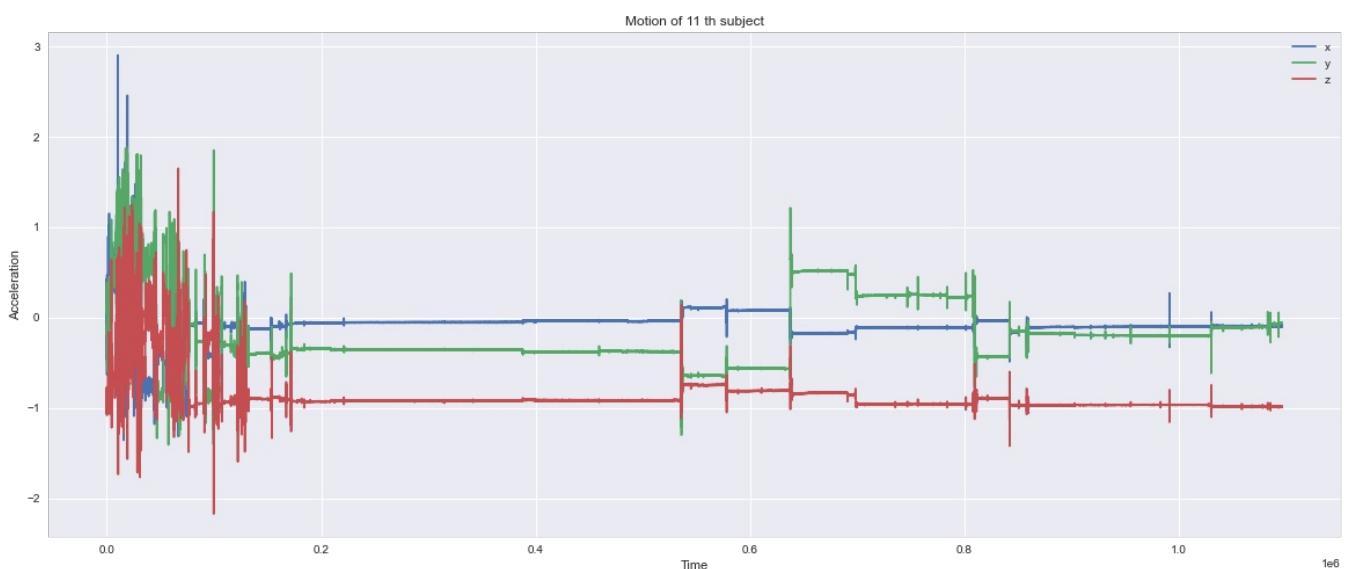
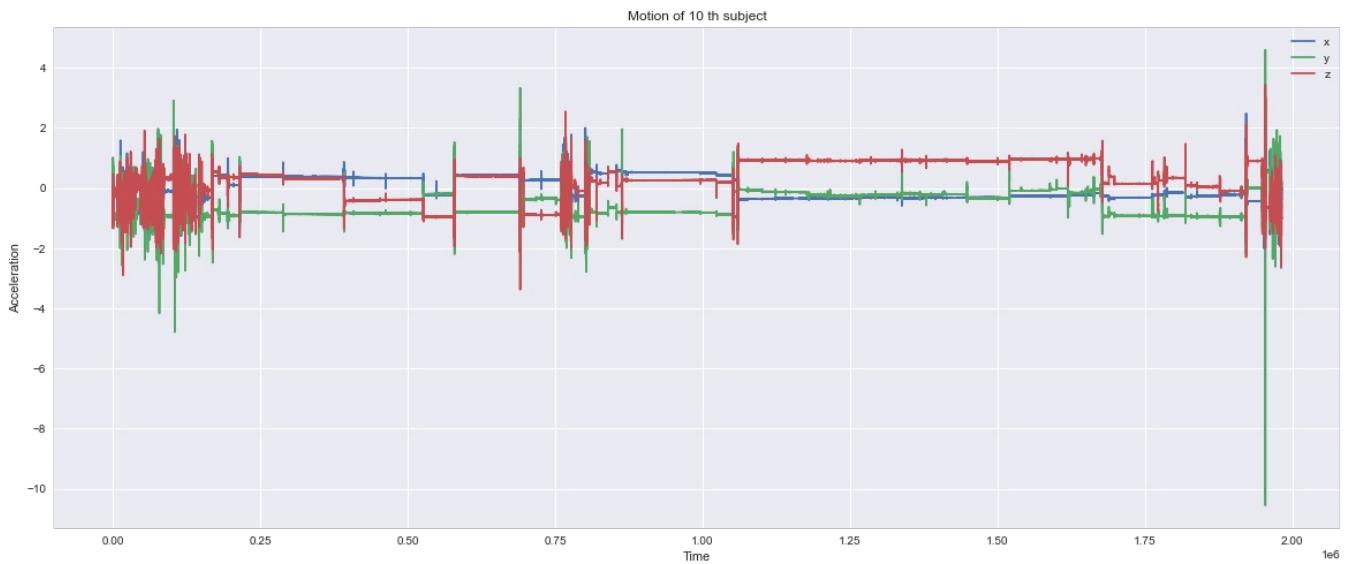
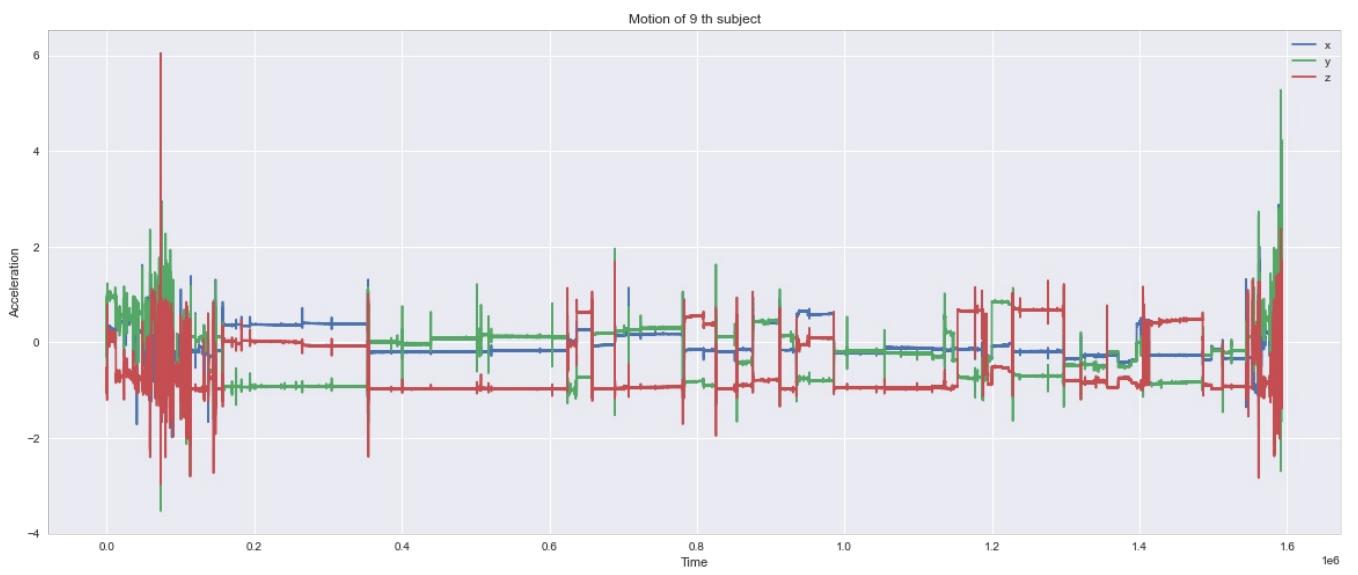
for i in range(len(motion)):
    plt.figure(figsize=(20, 8))
    plt.title(f'Motion of {i} th subject')
    plt.xlabel('Time')
    plt.ylabel('Acceleration')
    motion[i]['x'].plot(label='x')
    motion[i]['y'].plot(label='y')
    motion[i]['z'].plot(label='z')

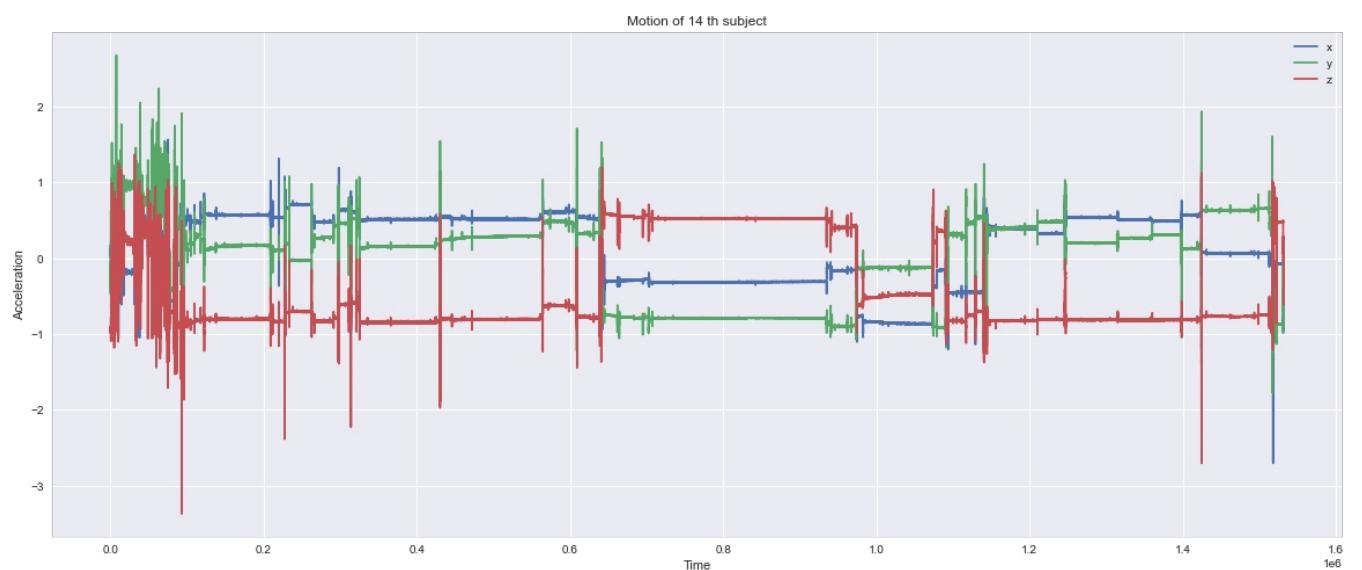
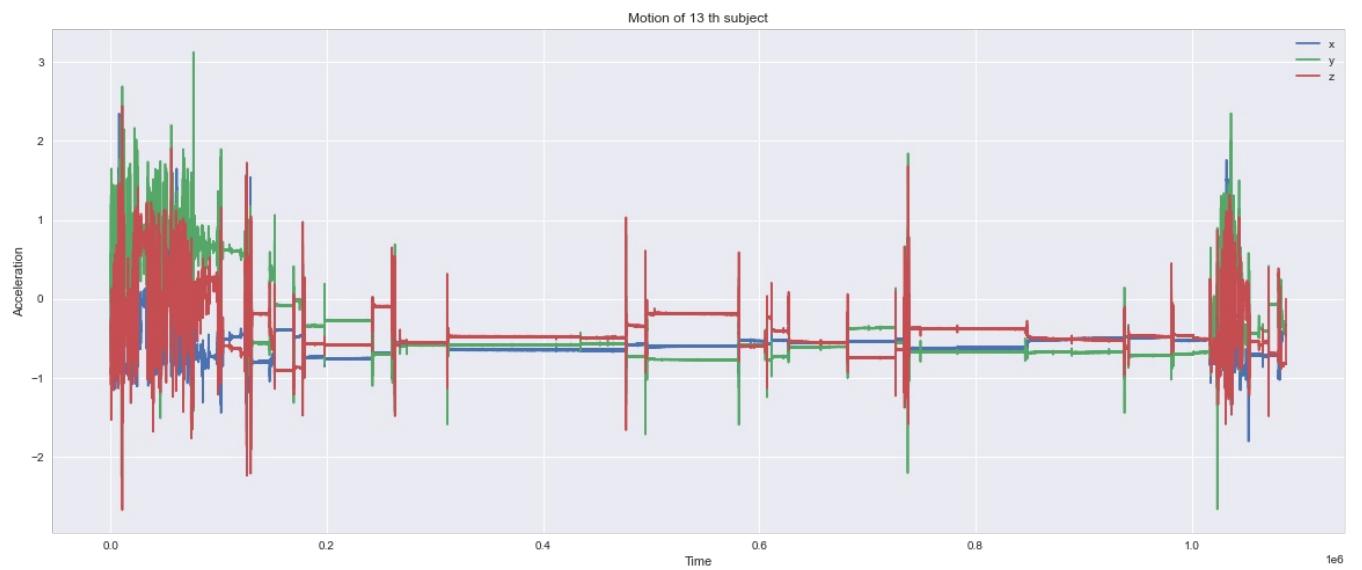
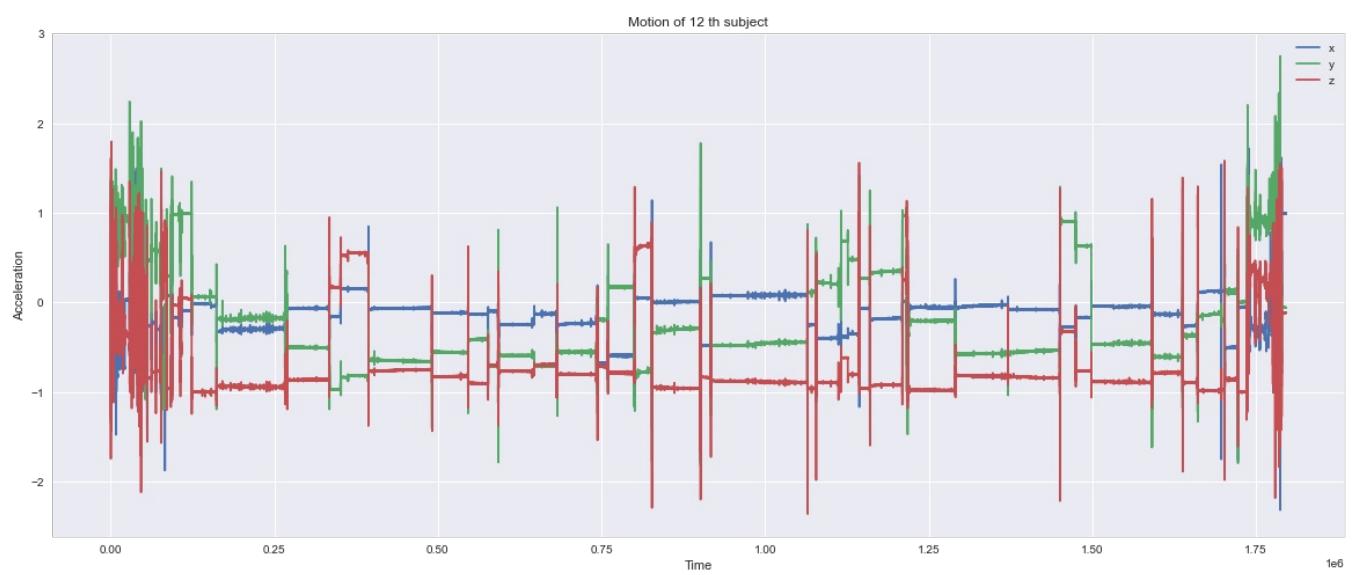
    plt.legend()
    plt.show()
```

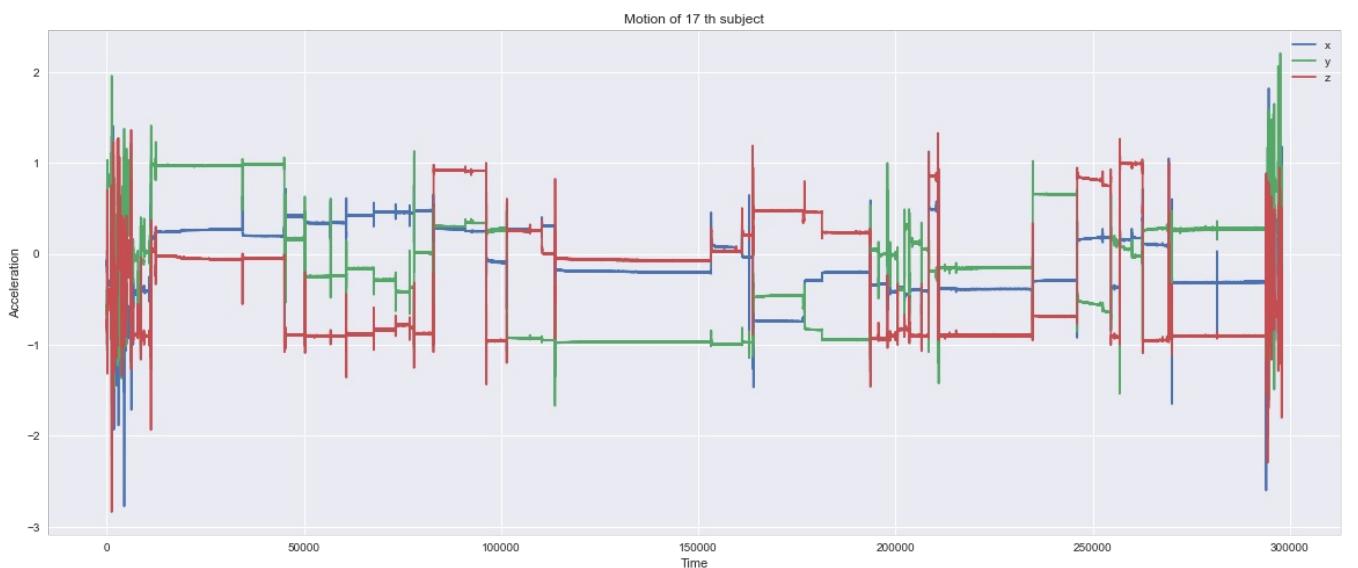
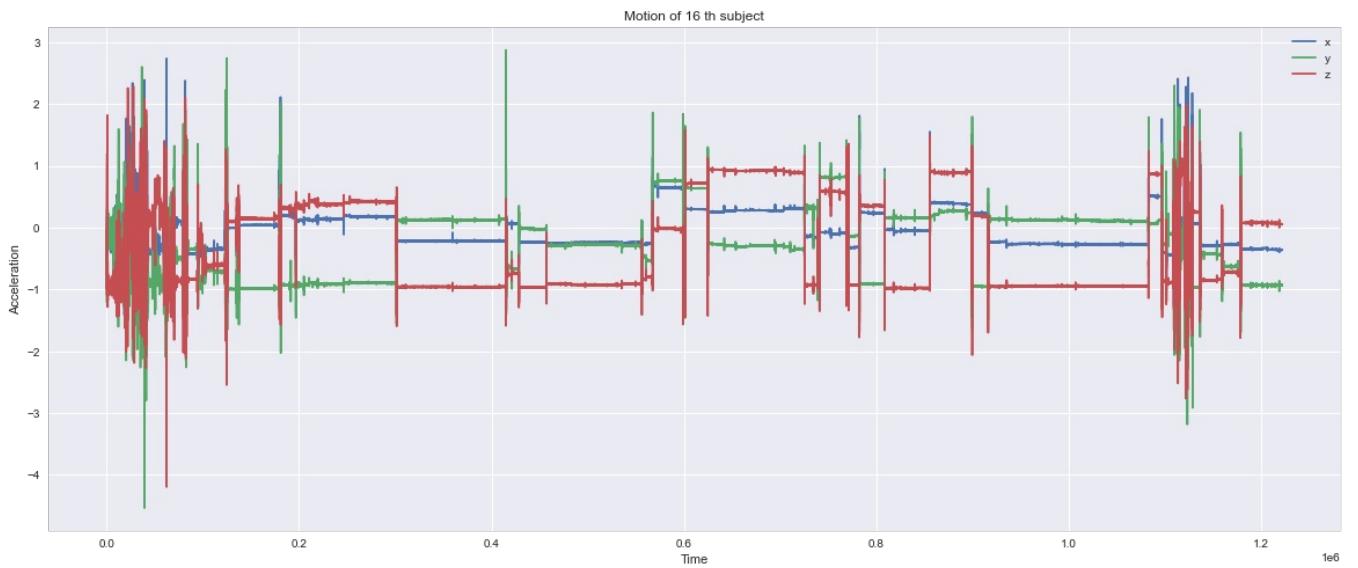
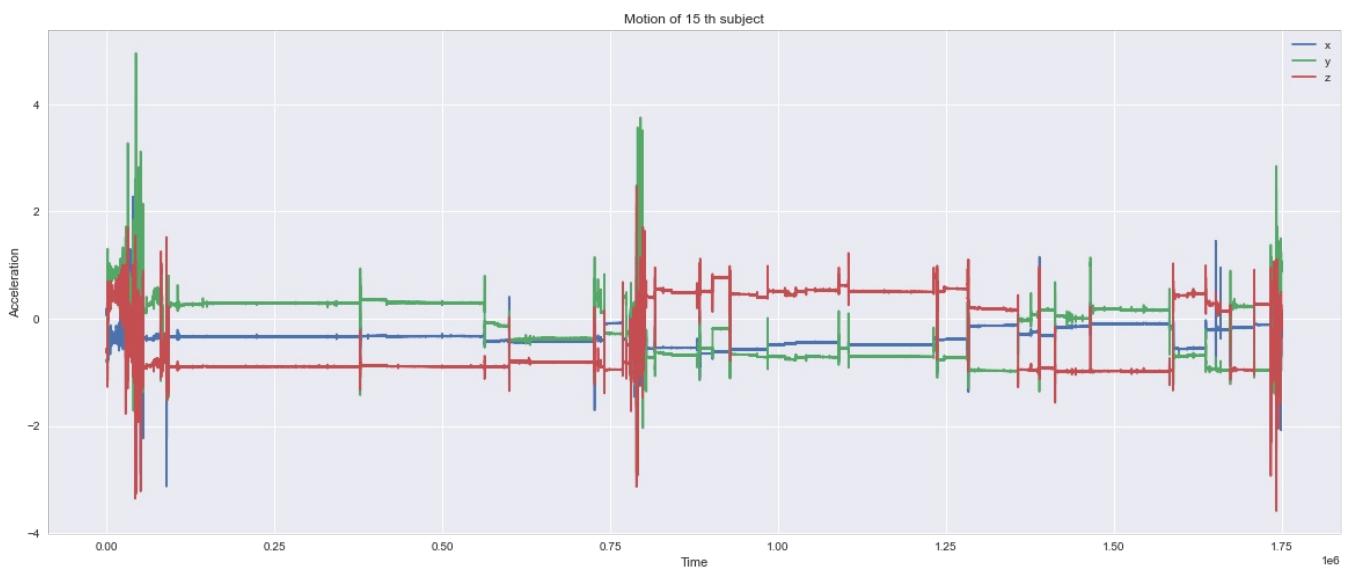


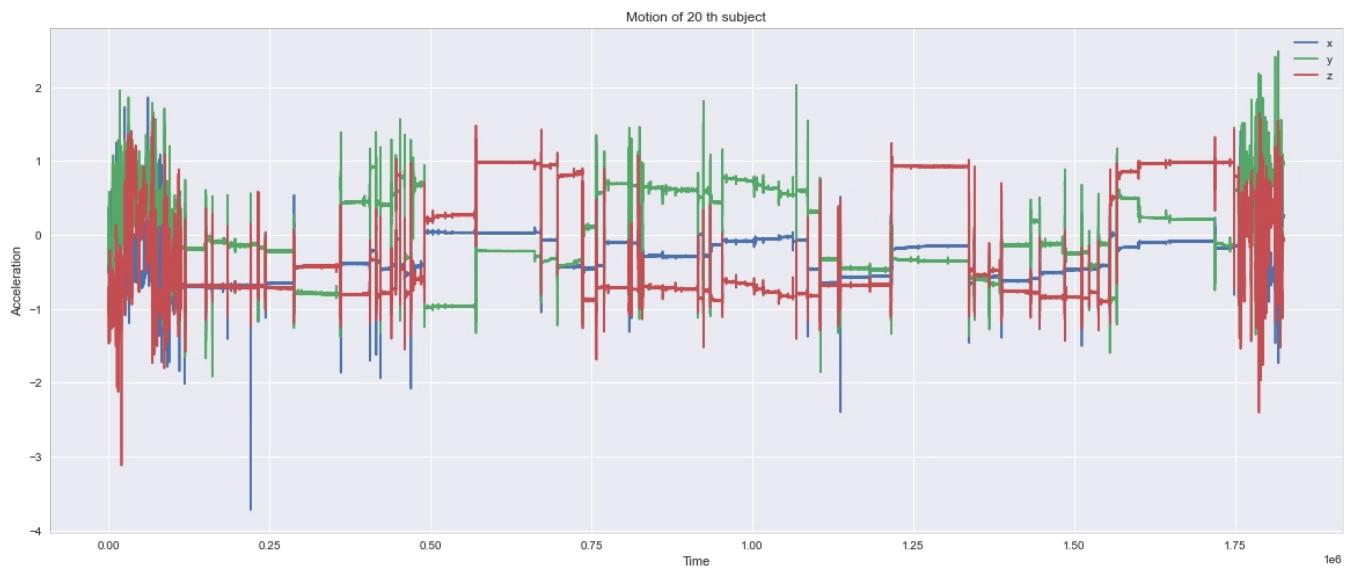
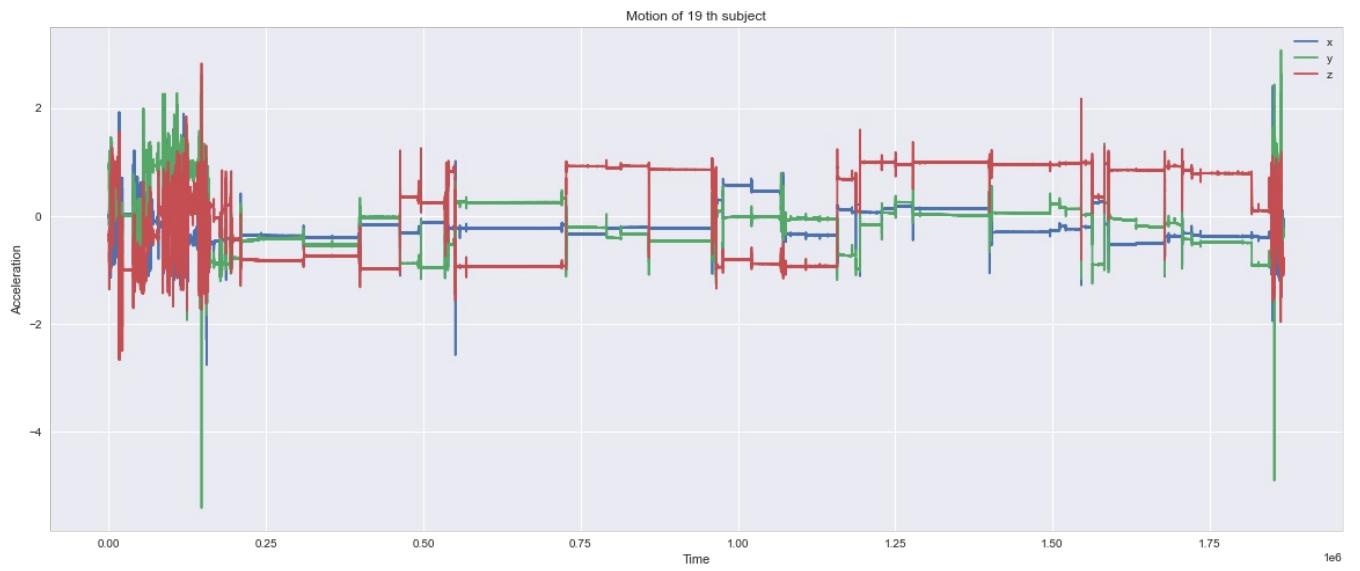
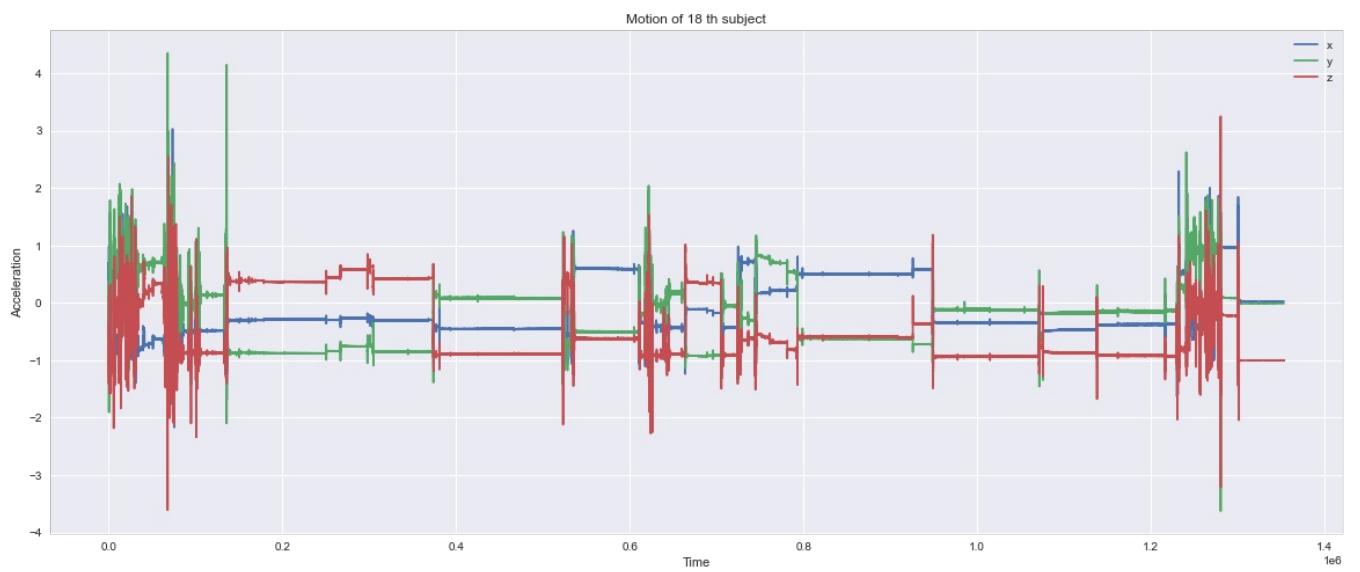


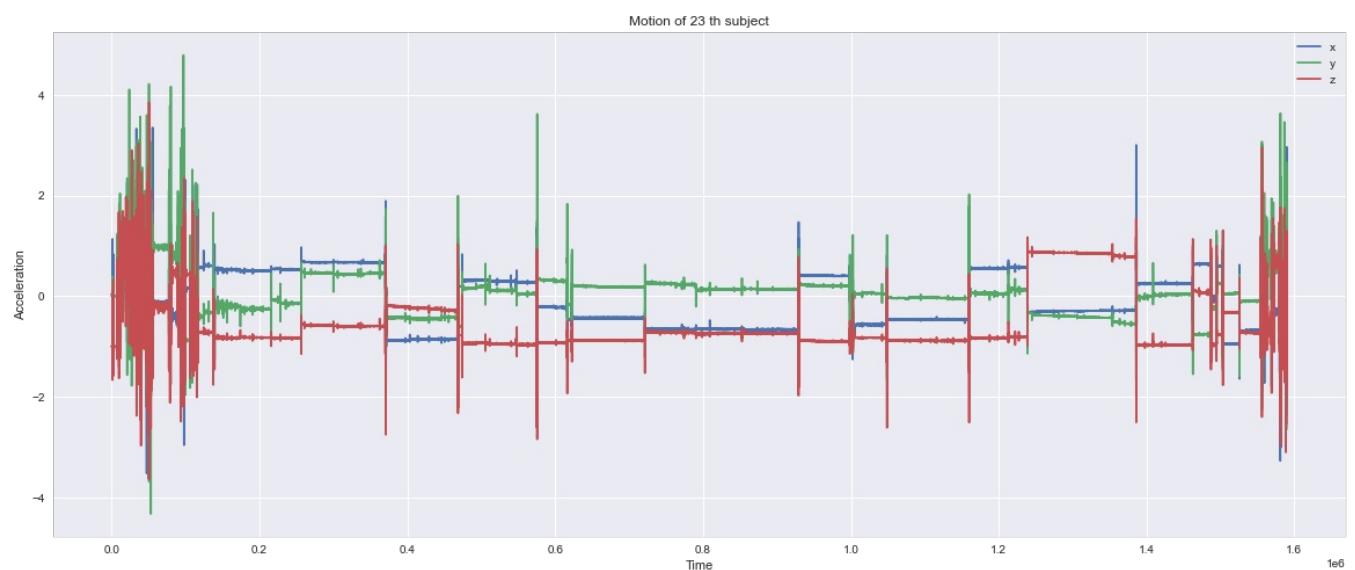
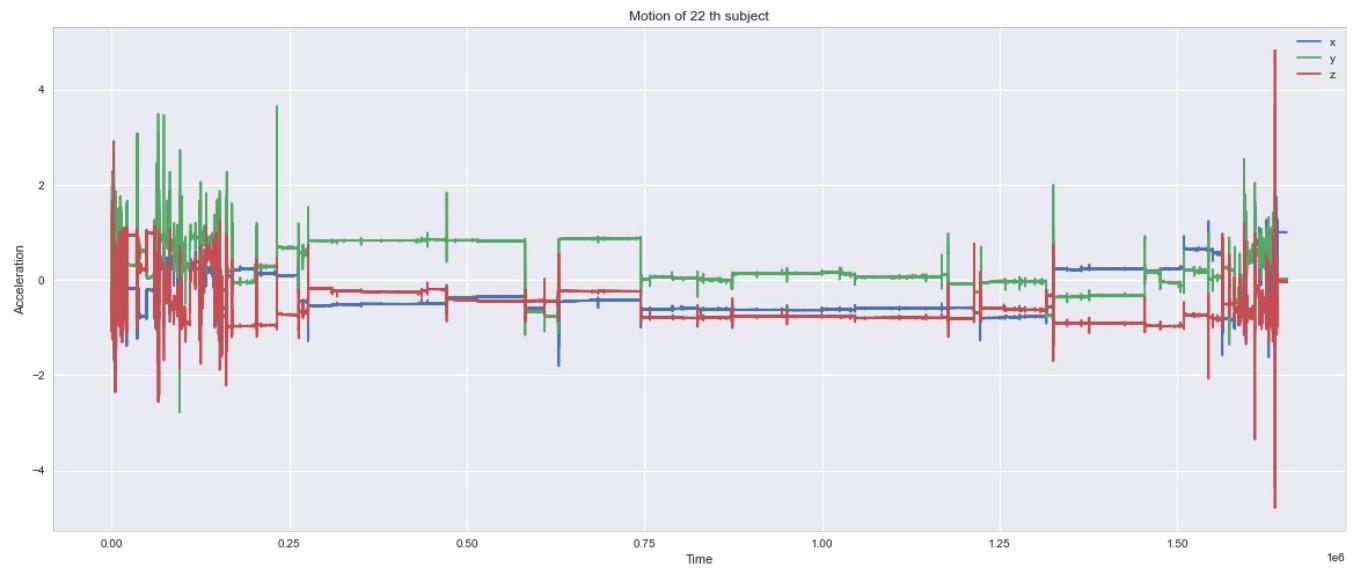
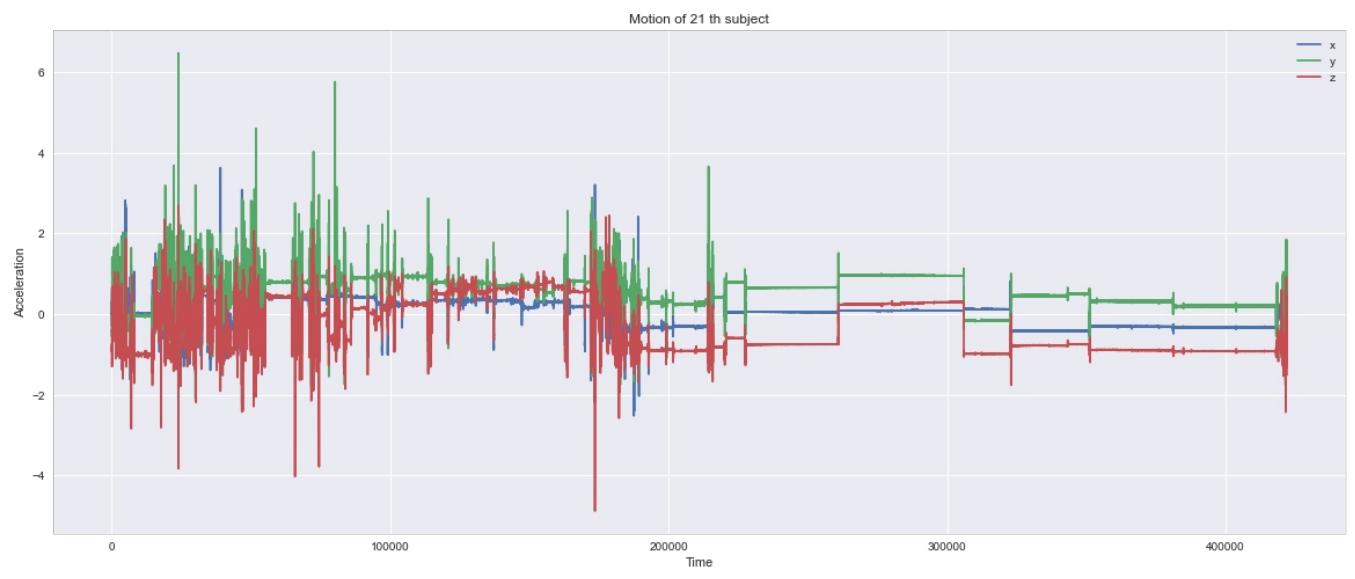


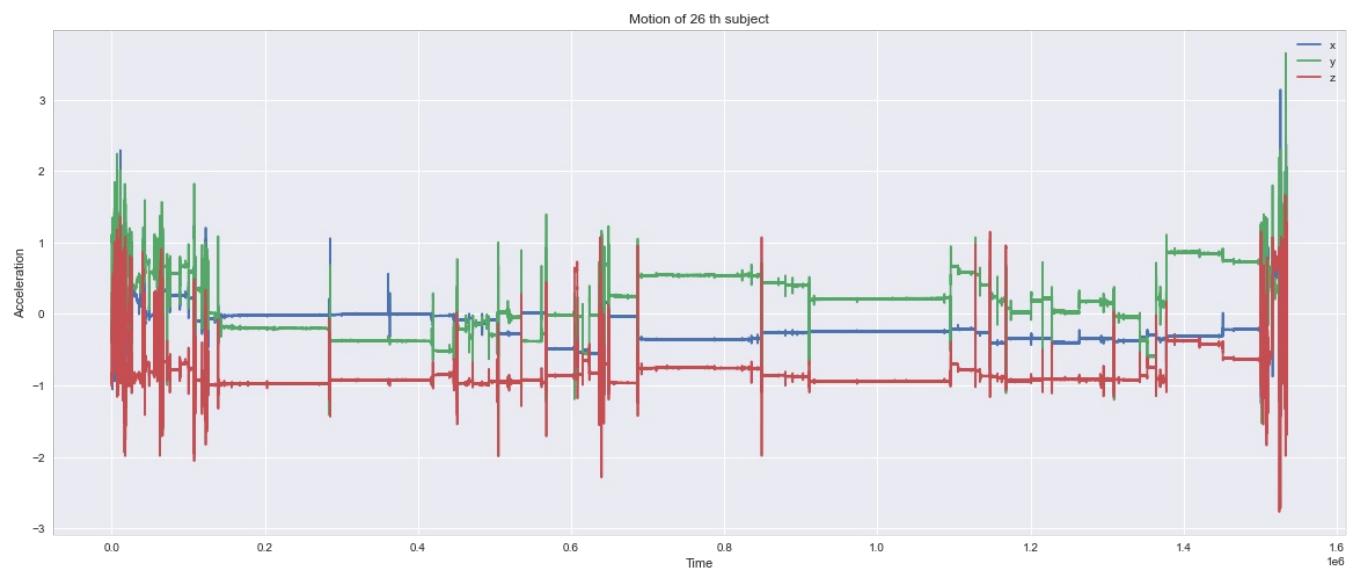
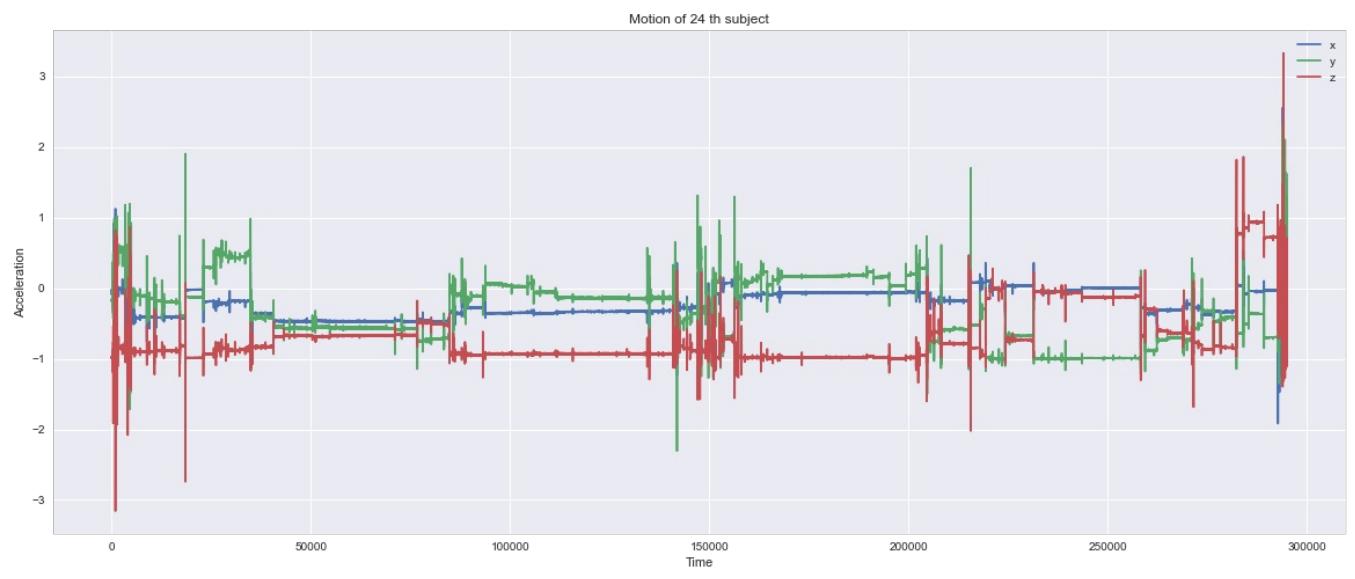


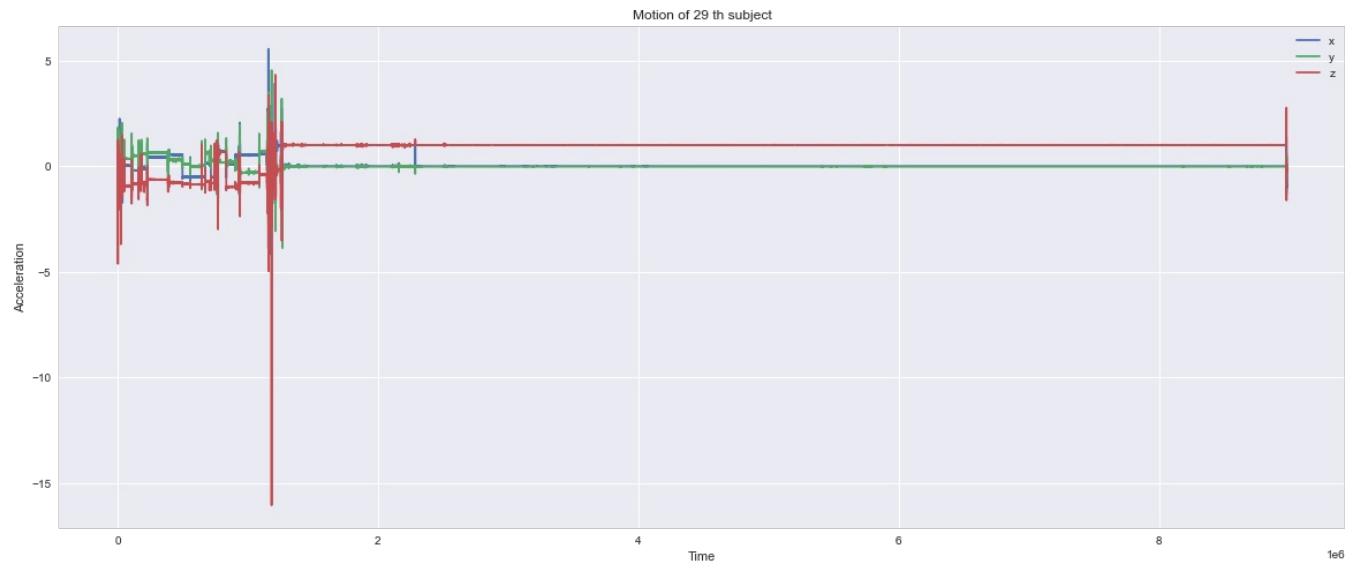
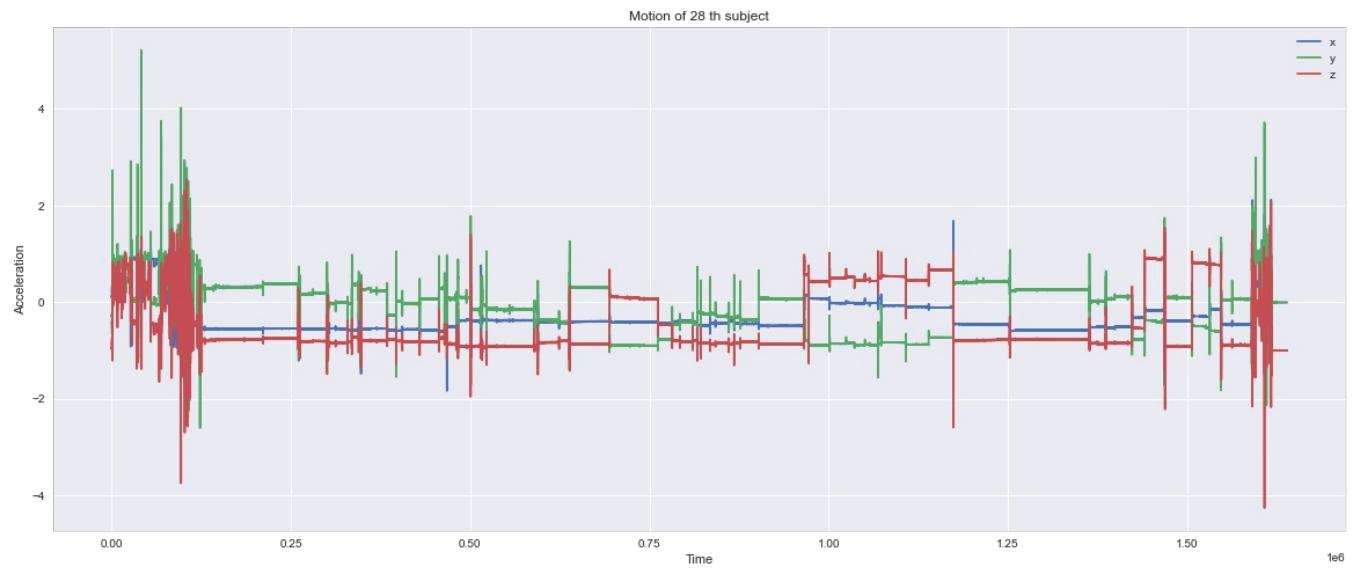
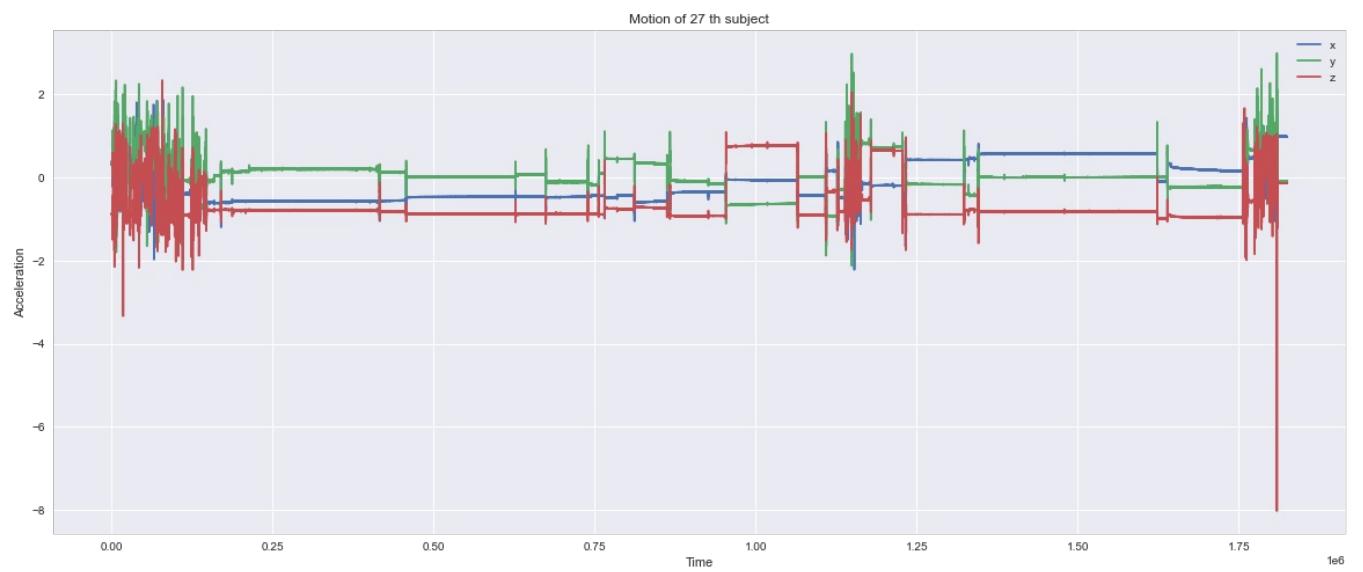


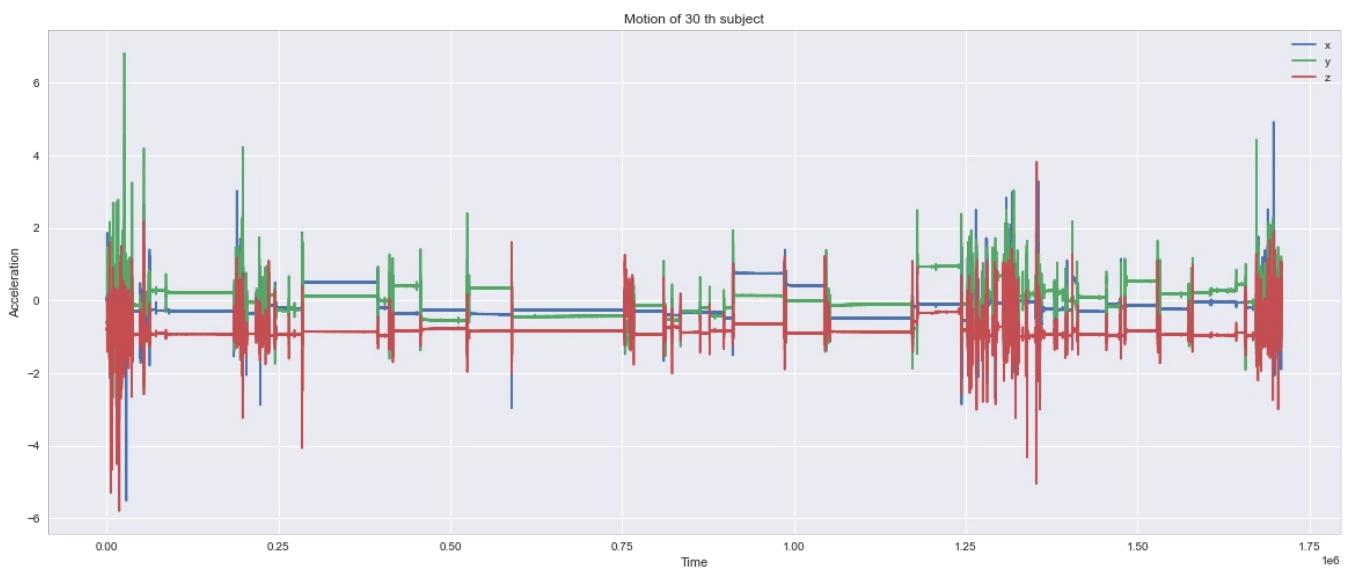












Label EDA

6, 17, 19, 22, 24, 25번 피험자는 중간에 -1(unscored)이 있으므로 수면 분류에 제외할 예정임.

17, 24번 피험자의 수면 라벨링에 설명에 나와있지 않은 4가 있으므로 수면 분류에 제외할 예정임.

```
In [16]: # label 데이터 프레임 리스트
label = []

dir_name = ''
path = ''
file_list = os.listdir(path)

# 주어진 파일 이름에서 숫자 부분을 추출
# 파일 이름에서 정수 값을 추출하여 파일을 숫자순으로 정렬하는 데 사용
def get_number(filename):
    return int(re.search(r'\d+', filename).group())

# 리스트 파일 이름 순으로 정렬
file_list_txt = sorted([file for file in file_list], key=get_number)

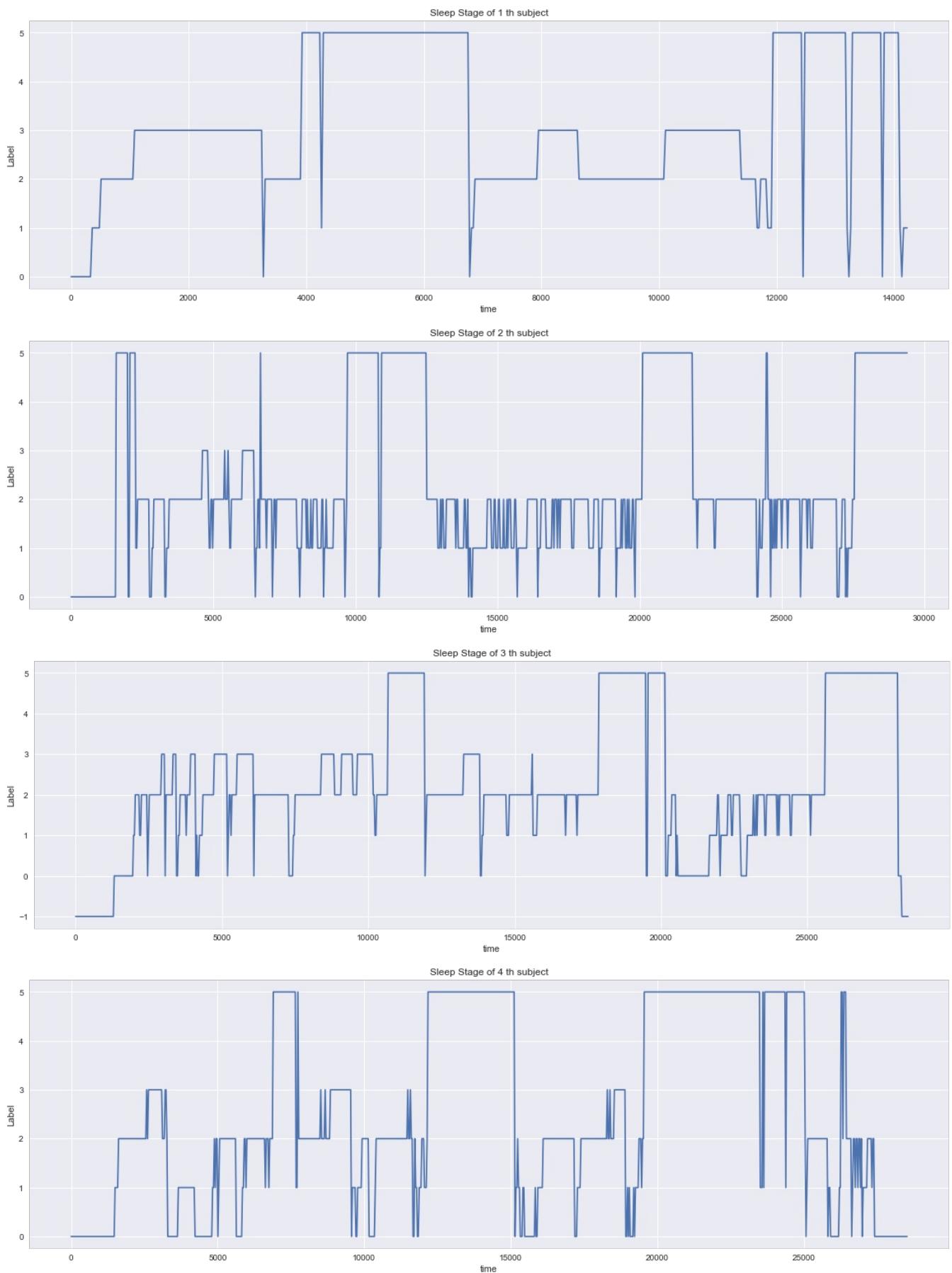
# 텍스트 파일을 csv 파일로 저장
for i in file_list_txt:
    df = pd.read_csv(dir_name+i, sep=' ', header=None, names=['time', 'label'])
    df = df.set_index('time') # 인덱스를 시간으로 지정
    label.append(df)
```

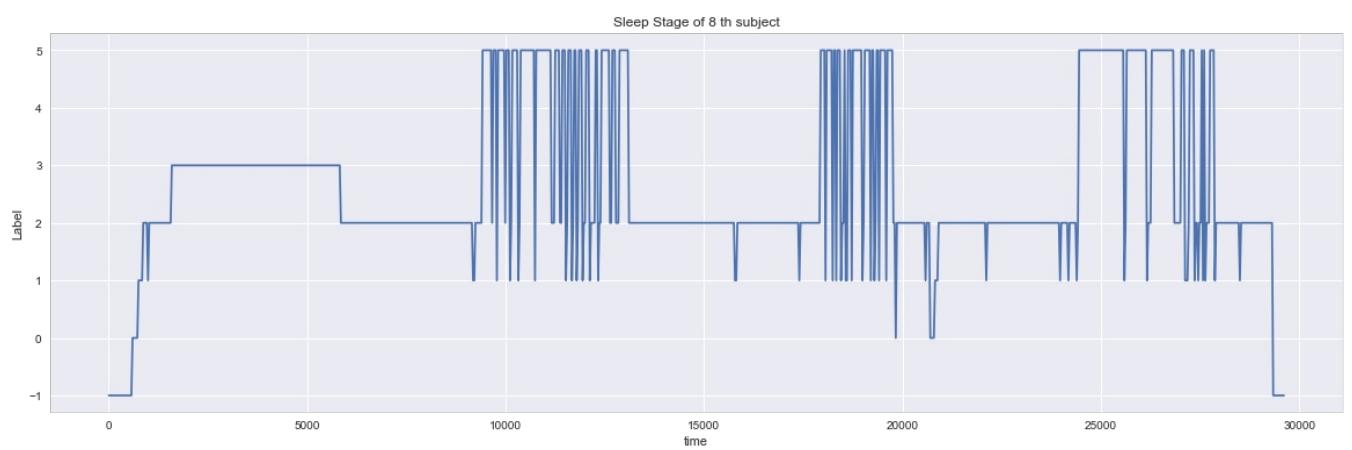
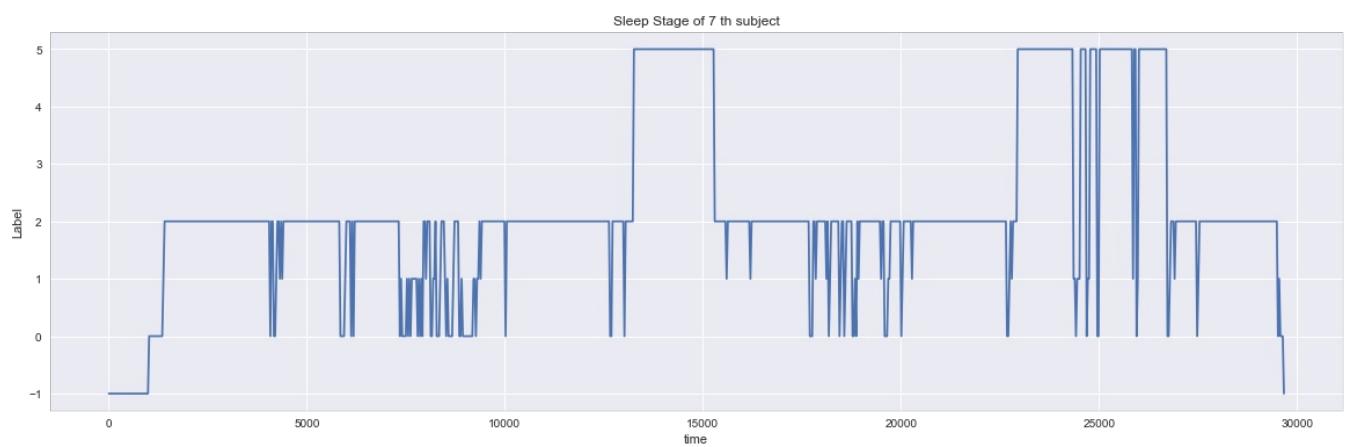
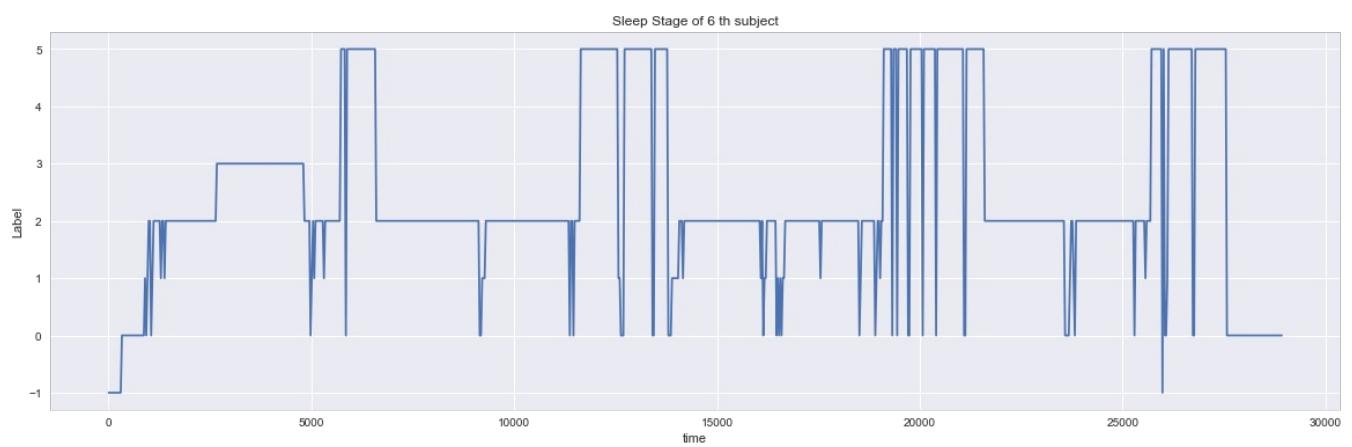
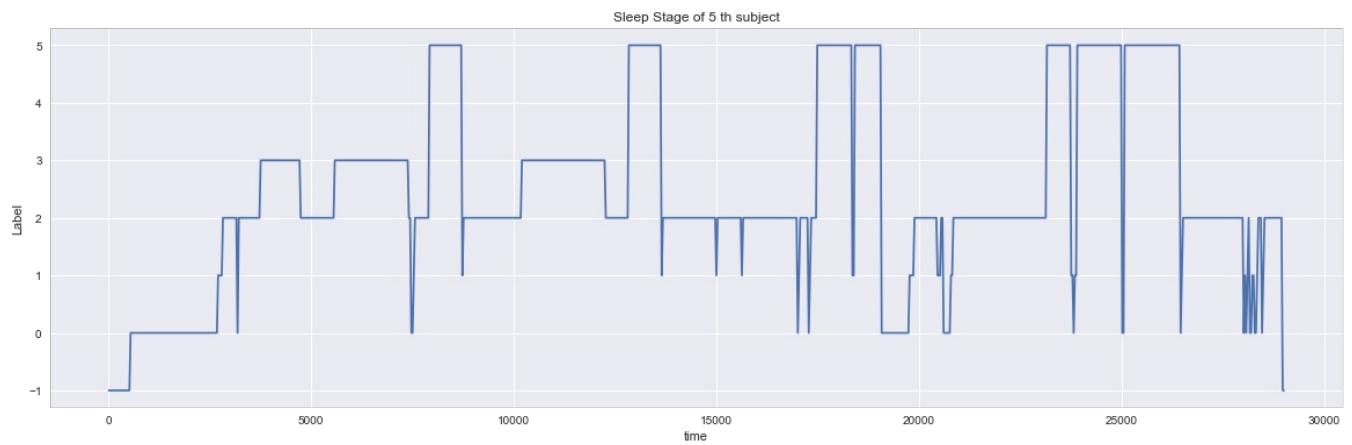
```
In [17]: # 6, 17, 19, 22, 24, 25 는 중간에 -1이 있으므로 제외

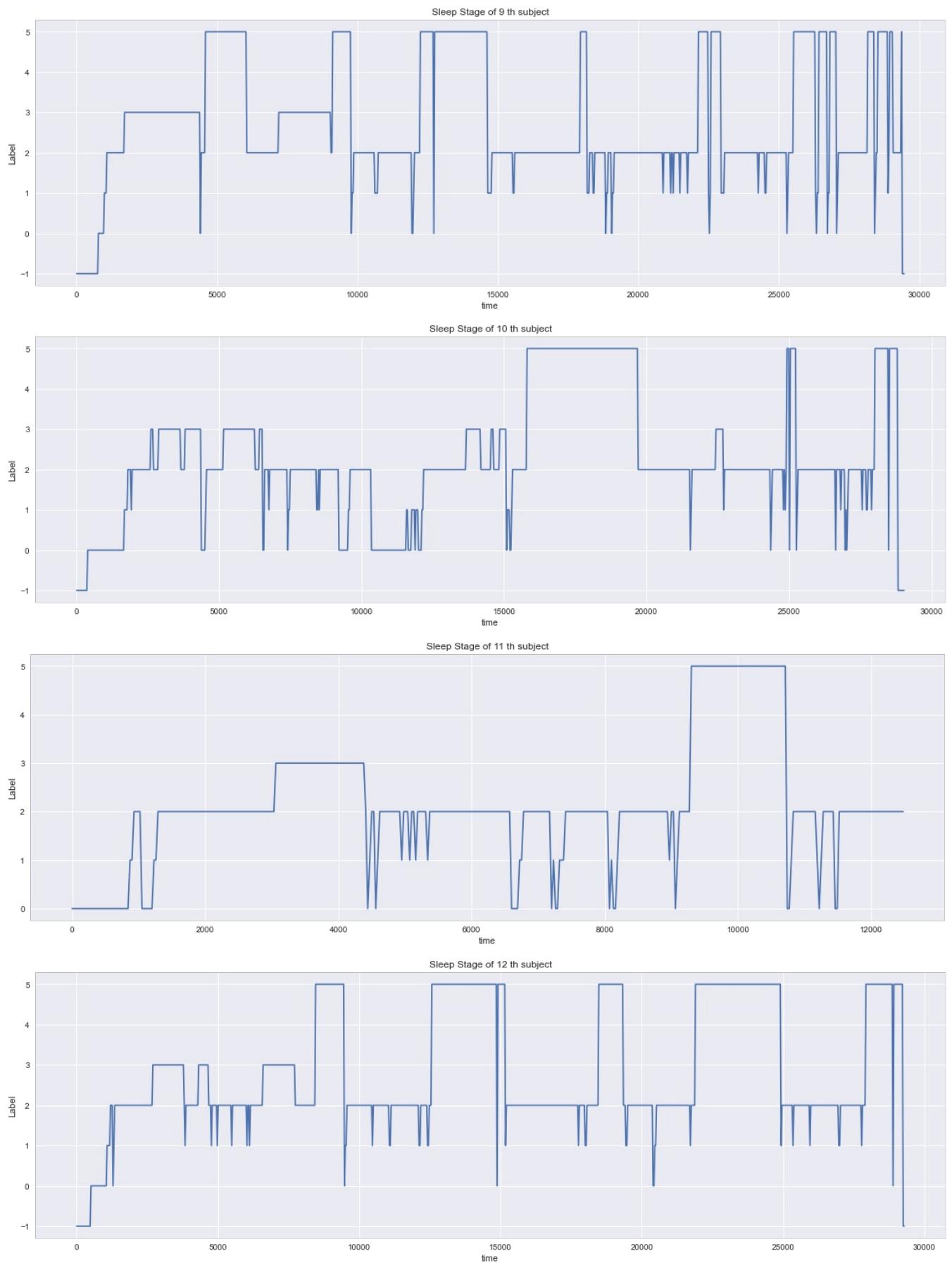
for i in range(len(label)):
    plt.figure(figsize=(20, 6))
    plt.title(f'Sleep Stage of {i} th subject')
    plt.xlabel('Time')
    plt.ylabel('Label')

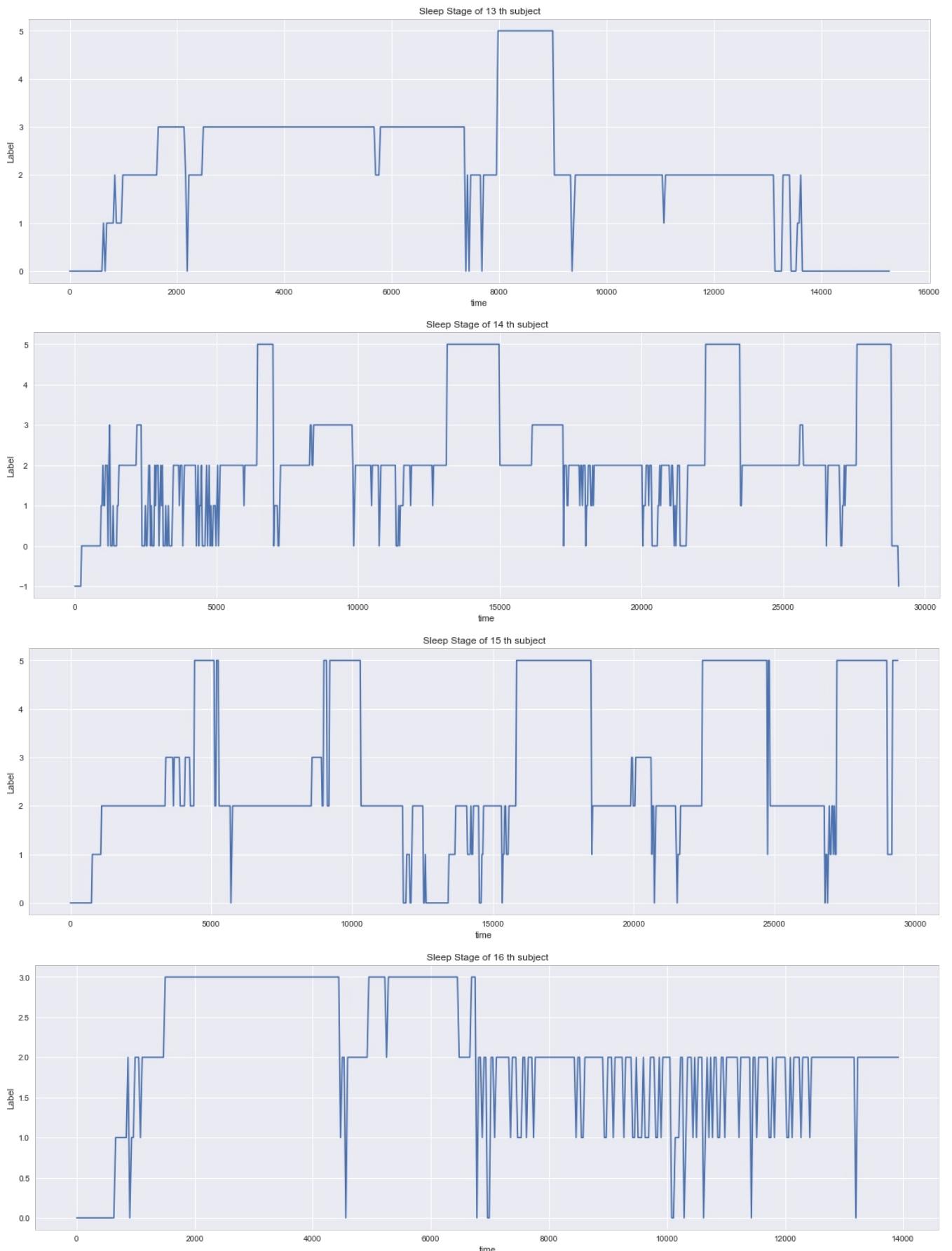
    label[i]['label'].plot()
    plt.show()
```

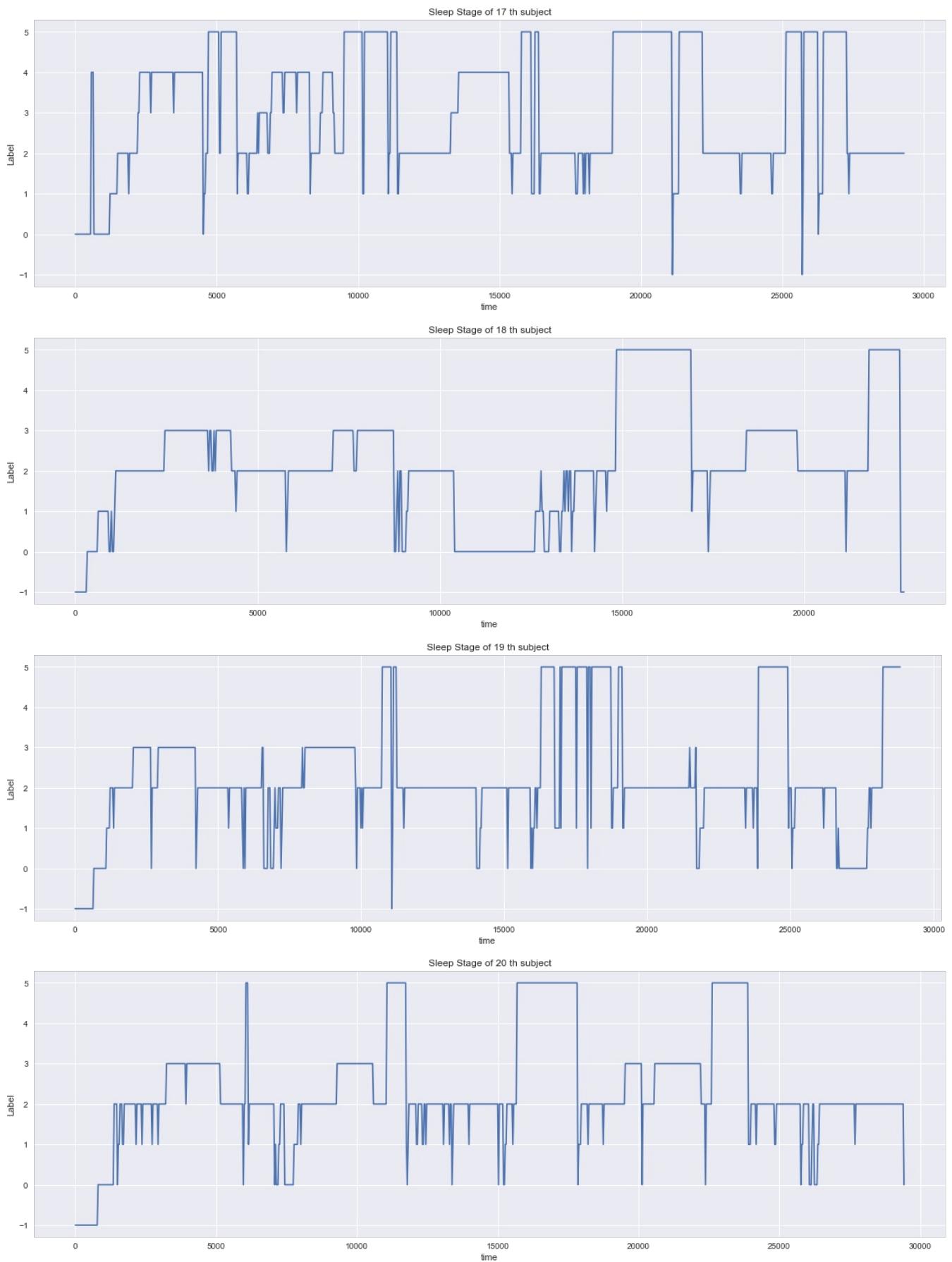


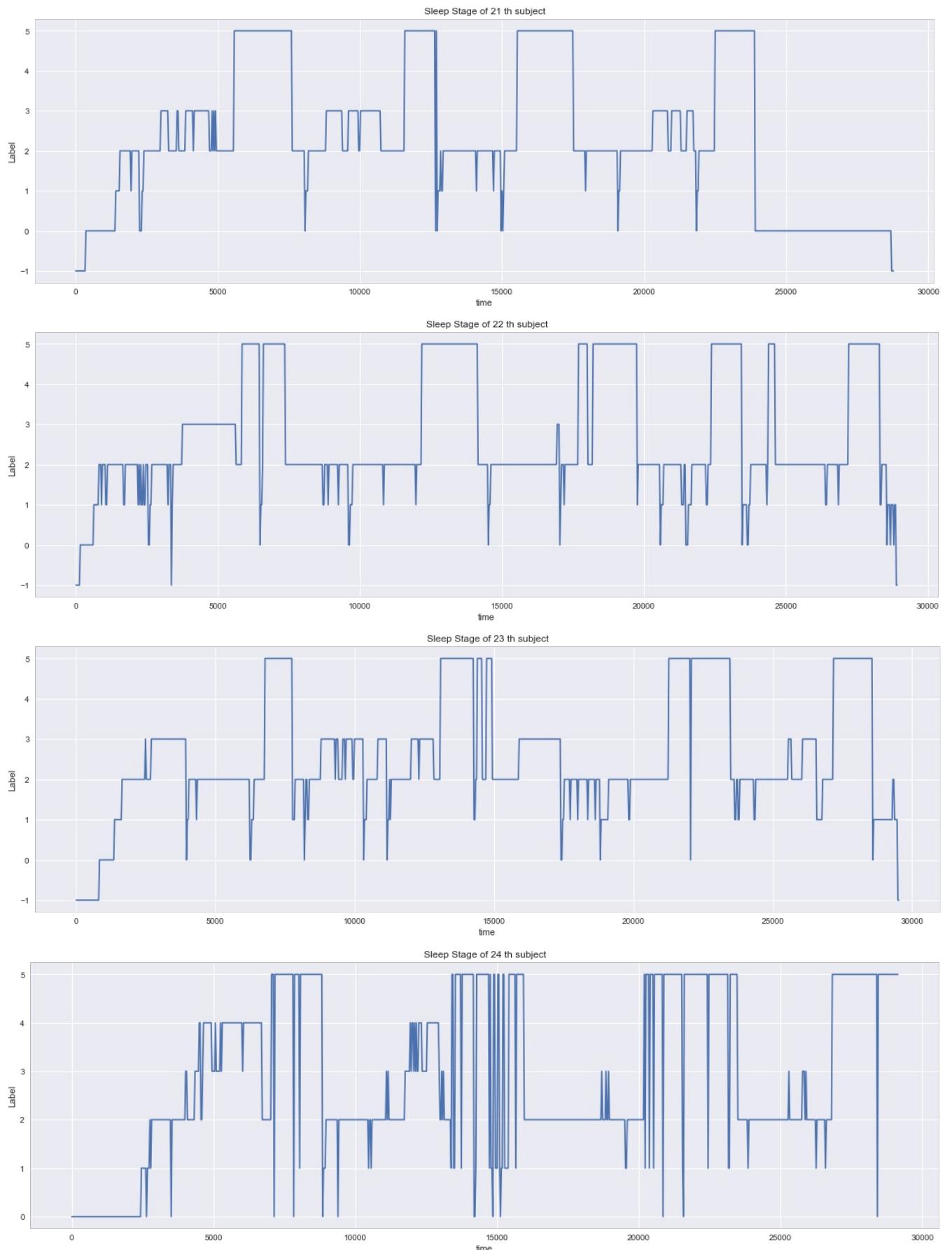


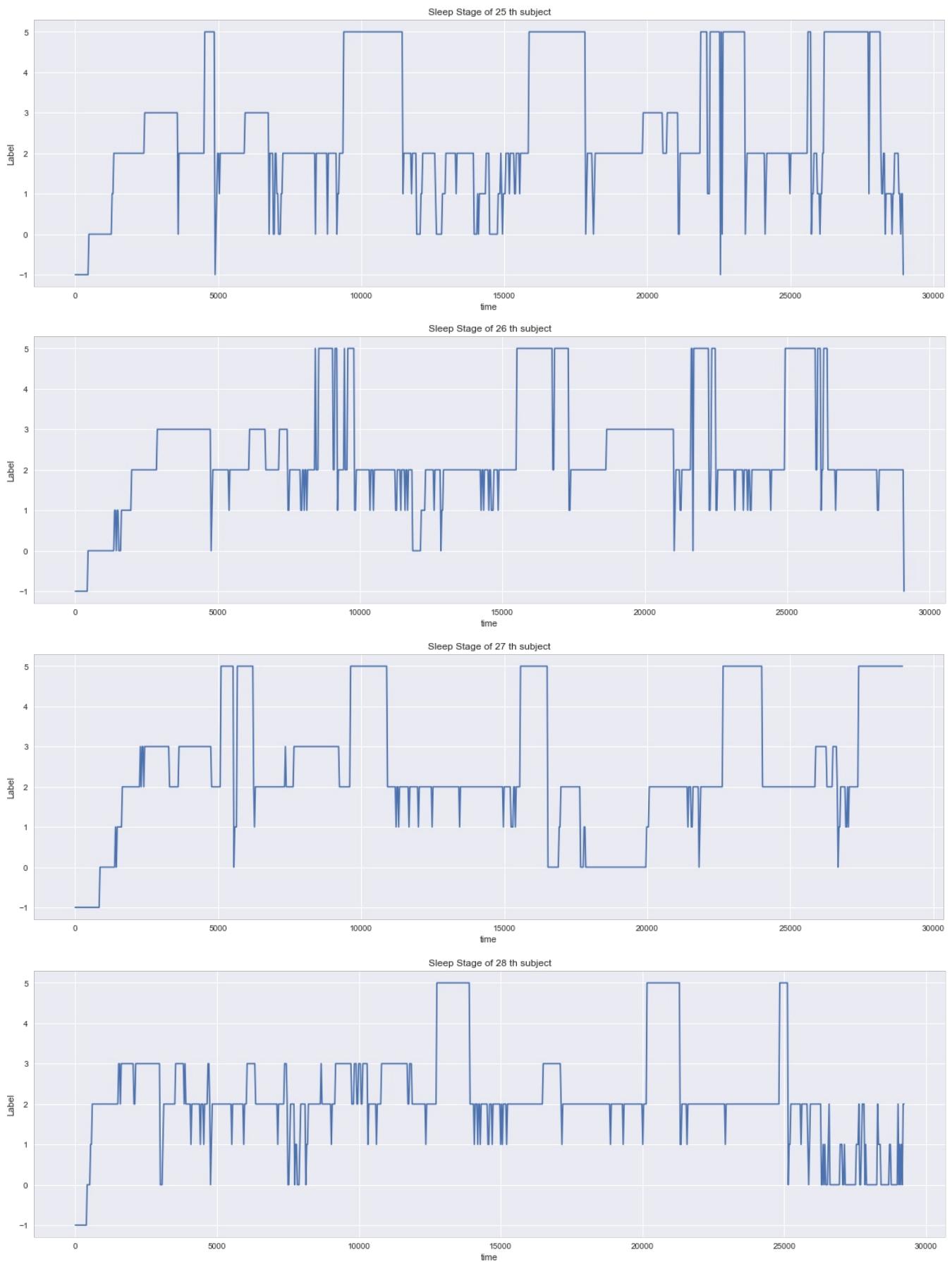


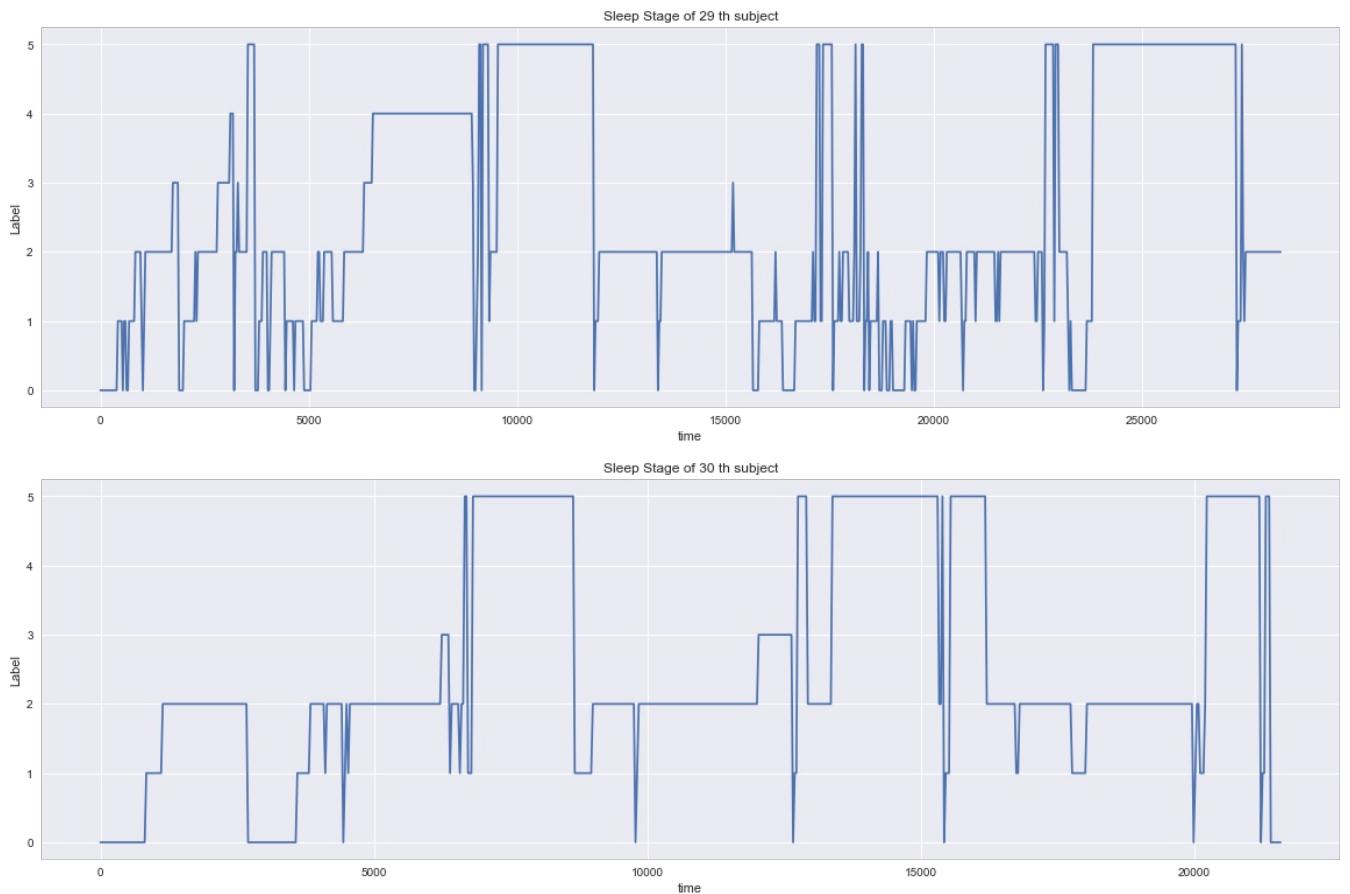




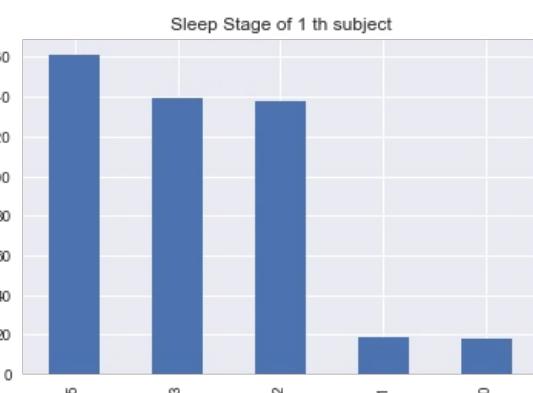
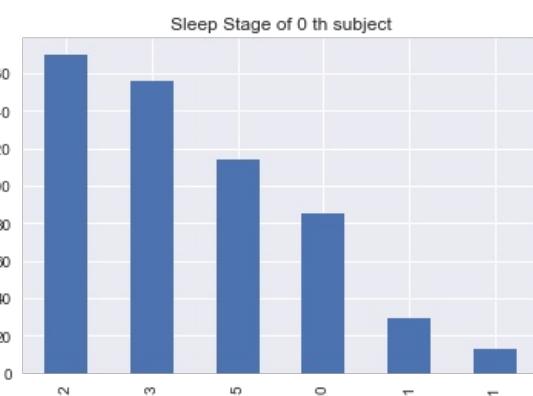




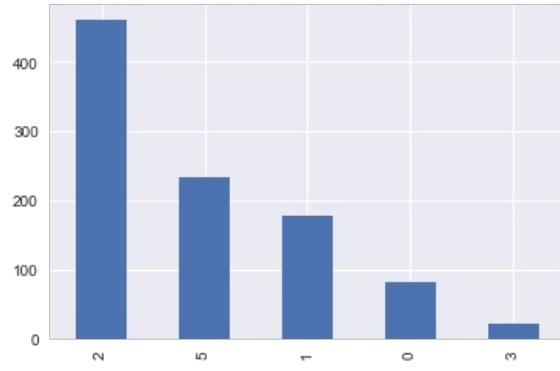




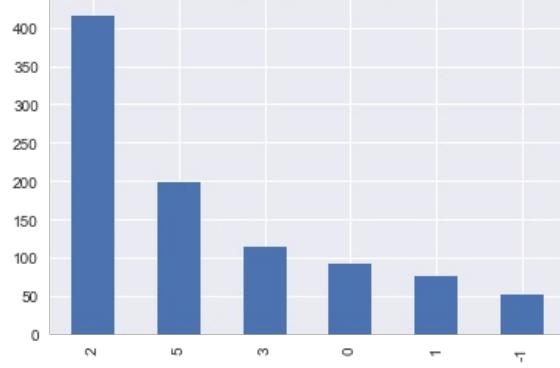
```
In [18]: # 각 수면 단계가 얼마나 있는지 확인
# 17, 24번 피험자의 수면 라벨링에 설명에 나와있지 않은 4가 있으므로 수면 분류에 제외
for i in range(len(label)):
    plt.title(f'Sleep Stage of {i} th subject')
    label[i]['label'].value_counts().plot.bar()
    plt.show()
```



Sleep Stage of 2 th subject



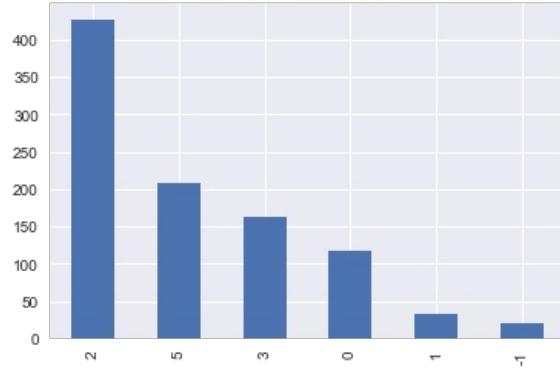
Sleep Stage of 3 th subject



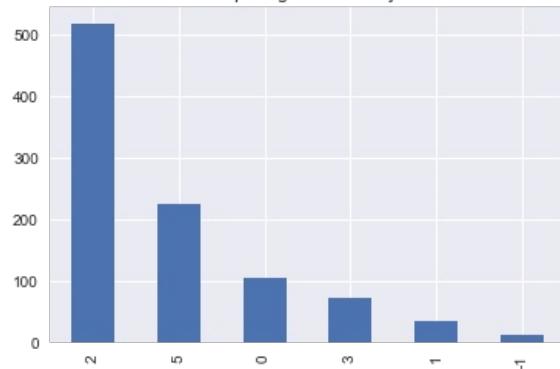
Sleep Stage of 4 th subject



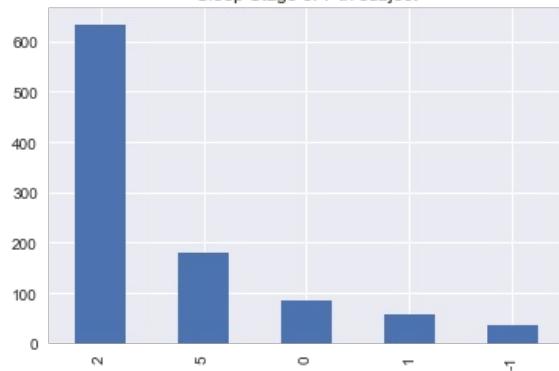
Sleep Stage of 5 th subject



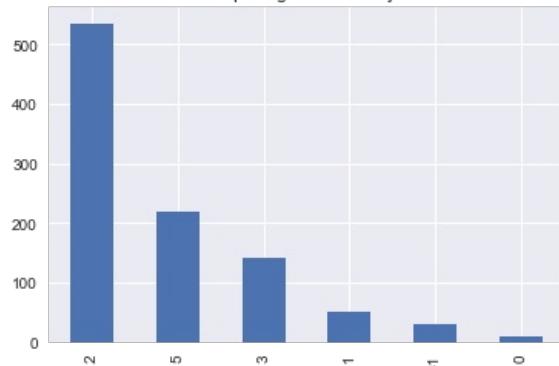
Sleep Stage of 6 th subject



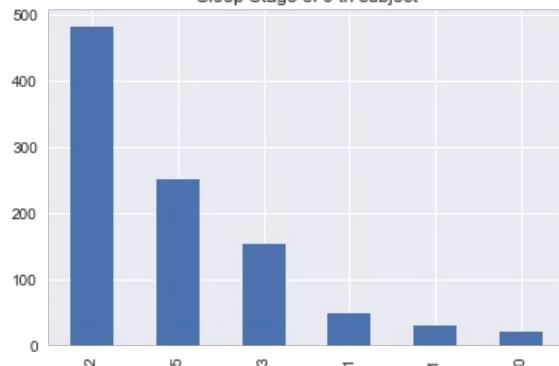
Sleep Stage of 7 th subject



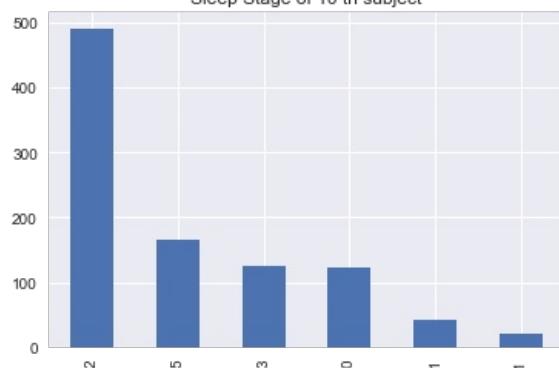
Sleep Stage of 8 th subject



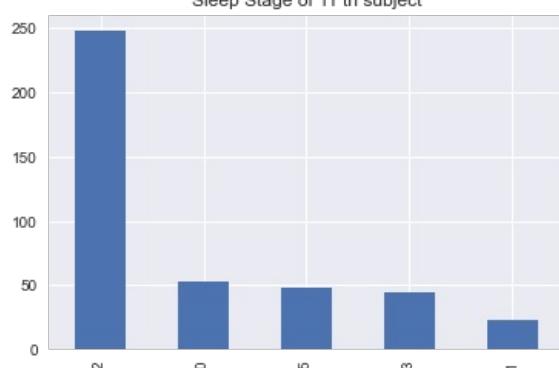
Sleep Stage of 9 th subject



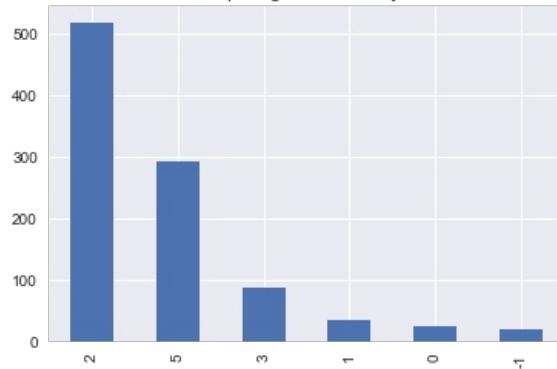
Sleep Stage of 10 th subject



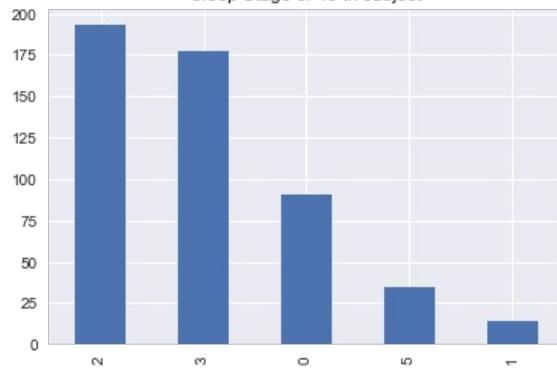
Sleep Stage of 11 th subject



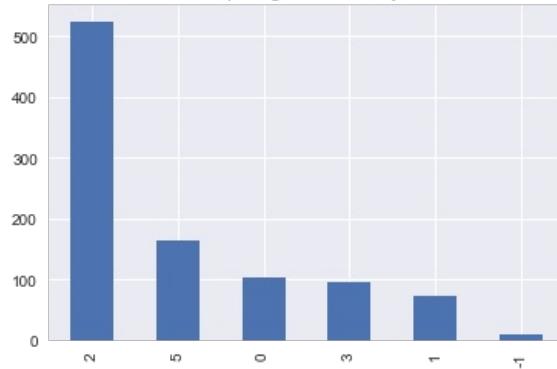
Sleep Stage of 12 th subject



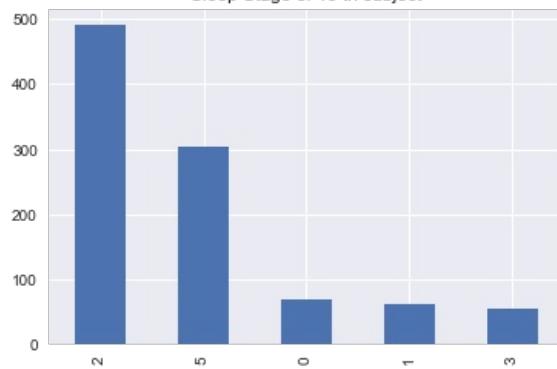
Sleep Stage of 13 th subject



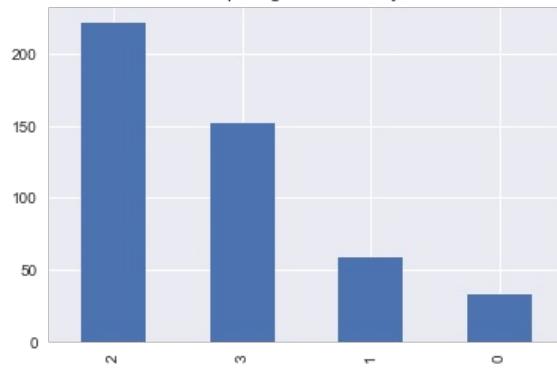
Sleep Stage of 14 th subject



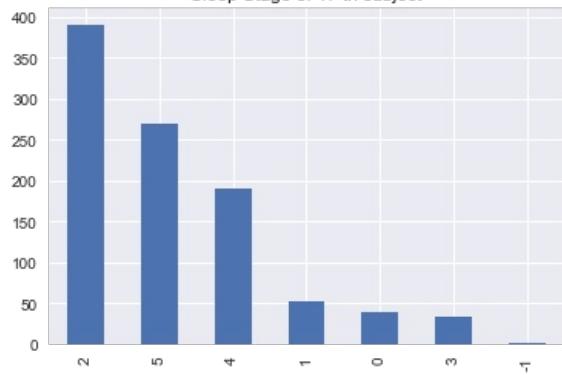
Sleep Stage of 15 th subject



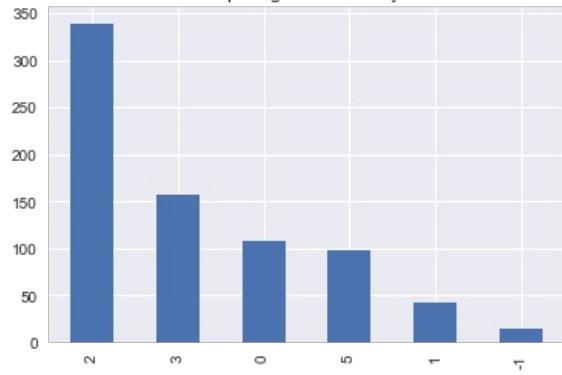
Sleep Stage of 16 th subject



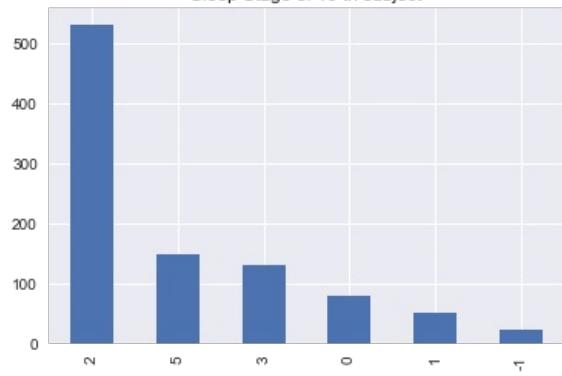
Sleep Stage of 17 th subject



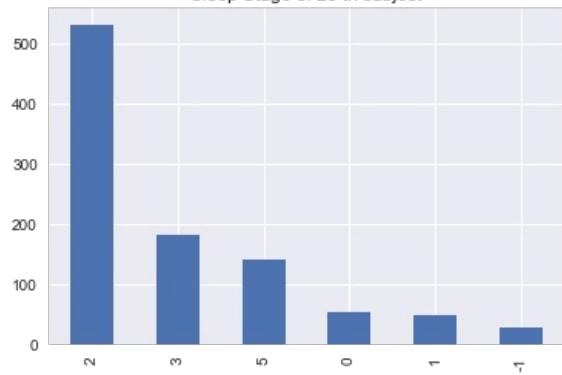
Sleep Stage of 18 th subject



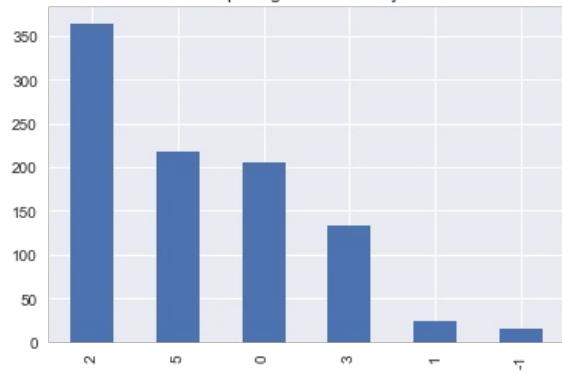
Sleep Stage of 19 th subject



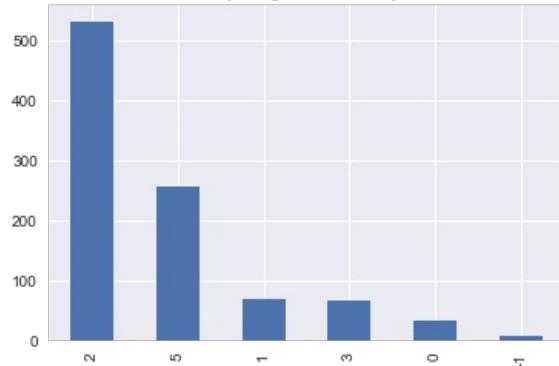
Sleep Stage of 20 th subject



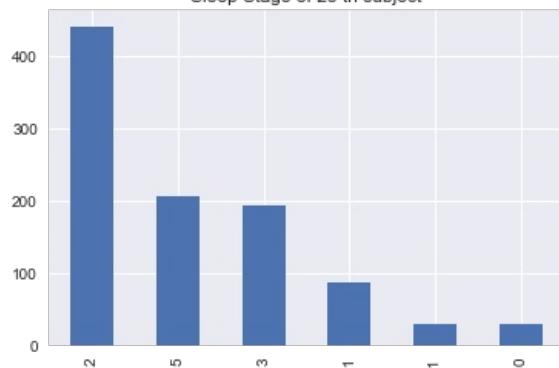
Sleep Stage of 21 th subject



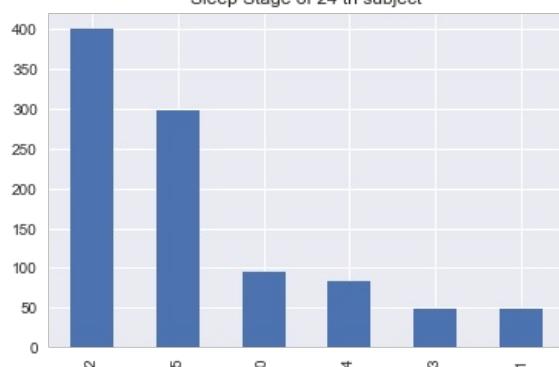
Sleep Stage of 22 th subject



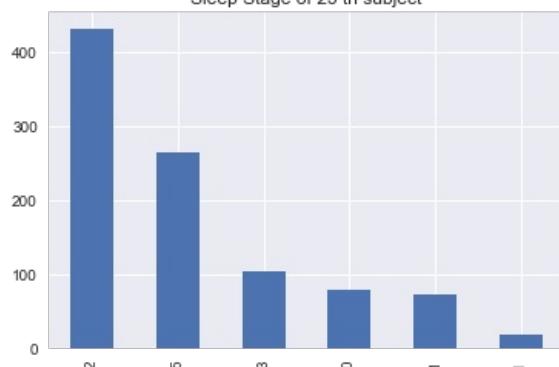
Sleep Stage of 23 th subject



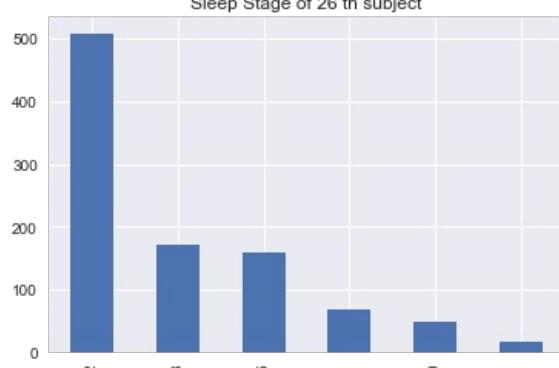
Sleep Stage of 24 th subject



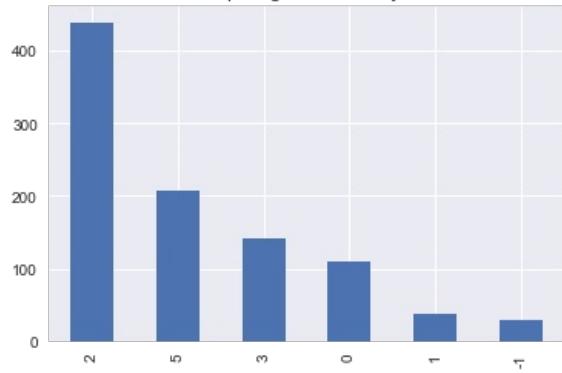
Sleep Stage of 25 th subject



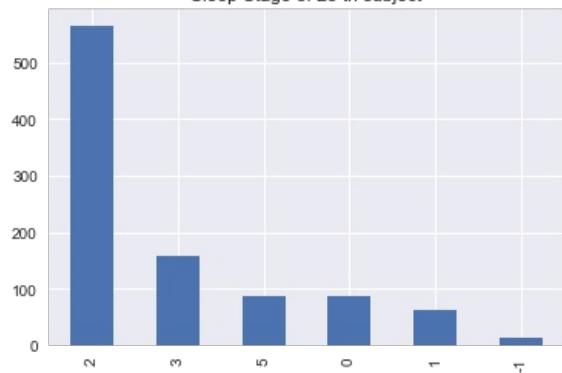
Sleep Stage of 26 th subject



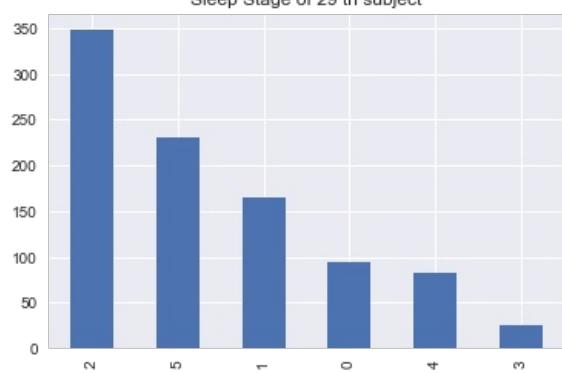
Sleep Stage of 27 th subject



Sleep Stage of 28 th subject



Sleep Stage of 29 th subject



Sleep Stage of 30 th subject

