

모바일 기기 도난 시 센서 데이터를 이용한 범죄자의 신장 추론

201920871 윤상훈

1. 주제 선정 이유

모바일 기기의 보급이 증가함에 따라, 이들 기기의 도난 사건도 빈번해지고 있습니다. 이러한 상황에서 도난된 기기를 신속히 추적하고 범인을 식별하는 것은 매우 중요합니다. 최근의 기술 발전으로 인해 모바일 기기는 다양한 센서를 탑재하여 사용자의 움직임과 행동을 추적할 수 있게 되었습니다. 기기가 움직일 때 수집되는 가속도, 중력, 회전 속도, 그리고 자세 데이터는 사용자의 체형과 성별을 정보를 추론하는 데 중요한 단서를 제공합니다.

본 프로젝트에서는 기기가 도난당했을 때 기기 내부 센서에서 수집된 데이터를 활용하여 범인의 체형, 그중 신장을 추론하는 방법을 탐구합니다. 제공된 데이터셋은 24명의 다양한 신체 정보를 가진 참가자들이 앉기, 걷기, 계단 오르기 및 내리기 등의 일상적인 동작을 수행하는 동안 수집된 센서 데이터를 포함합니다. 위 데이터 중 걷기와 달리기(조깅)에 대한 데이터를 분석하여 개인의 신체적 특징(신장)을 정확하게 식별할 수 있는지 평가하고, 도난 사건에서 잠재적인 범인을 추적하는 데 활용할 수 있는지를 조사합니다.

하지만 본 주제에는 몇 가지 한계가 존재합니다. 먼저, 기기를 소지하고 있는 사람의 정보를 특정하기에는 표본 크기가 24명으로 한정되어 있어 사용자의 신장을 추론하는 데 한계가 있을 수 있습니다. 또한, 이러한 기술이 오히려 범죄에 악용될 수 있다는 윤리적 문제도 제기될 수 있습니다. 그러나 다양한 센서 데이터를 통합하여 유의미한 결과를 도출할 것이며, 본 핸드폰의 주인만이 이러한 기능을 사용할 수 있도록 하는 등 적절한 방법을 통해 악용에 대한 문제를 해결할 수 있을 것입니다.

이러한 데이터 분석을 통해 도난 사건에서 범죄자 추적 및 식별 방법을 개선하고, 향후 보안 시스템의 발전에 기여할 수 있는 가능성을 모색합니다. 더 나아가, 더 큰 규모의 데이터를 수집하고 다양한 변수를 추가함으로써 이러한 모델을 더욱 발전시킬 수 있을 것입니다.

2. 수행 과정 개요

데이터셋에는 24개의 데이터 주체에 대해 6가지 동작(계단 오르기, 계단 내리기, 앉기, 서기, 걷기, 조깅하기)에 대한 센서 데이터로 구성되어있다. 본 프로젝트에는 걷기와 조깅 동작이 가장 핵심적인 동작이 될 것이다. 수집된 센서 데이터를 통해 걷기 동작인지, 조깅 동작인지를 구분하고 동작에 따른 신장을 구분할 것이다.

이와 같은 방법이 가능한 이유는 일반적으로 보폭은 신장에 비례하기 때문이다¹⁾. 이번 프로젝트에서는 구체적인 보폭의 길이를 직접 구하는 것이 아닌 서로 다른 보폭에서 발생하는 센서 데이터의 특징에 주목하여 신장을 예측할 것이다. 따라서, 예측 모델은 신장에 따라

1) <https://www.verywellfit.com/set-pedometer-better-accuracy-3432895>

위 문서에서는 여성의 경우 신장(inches)에 0.413, 남성의 경우 키에 0.415를 곱하면 대략적인 보폭을 얻을 수 있다고 한다.

비례하는 선형적인 모습을 관측할 수 있을 것으로 예상된다.

본격적인 분석에 앞서 수월한 분석을 위해 데이터 전처리 혹은 시각화하는 함수 등을 만들고 4개의 서로 다른 타입의 데이터(자세, 중력, 가속도, 회전) 데이터에 대한 비교 분석을 하여 적절한 데이터를 택할 것이다. 다만, 일반적으로 핸드폰 센서를 통해 걸음 수를 분석할 때, 가속도 센서를 사용함으로 가속도 데이터를 사용하는 것은 매우 권장될 것이다. 이러한 점에 유의하여 신장을 구분하기 위한 적절한 데이터를 택할 것이다.

전반적인 학습 과정에서 두 객체에 대한 데이터를 검증을 위해 분리할 것이다.²⁾ 모델 학습에 대한 모든 과정을 마치고 해당 데이터를 사용하여 결과를 확인하여 모델에 대한 최종 검증을 마칠 것이다.

3. 수행 과정

#데이터 로드

데이터를 가져오는 과정은 깃허브 공식 문서³⁾의 주어진 코드를 수정하여 구성하였다.

get_ds_infos 함수	creat_time_series 함수
<pre> 1 def get_ds_infos(): 2 """ 15 16 dss = pd.read_csv("data_subjects_info.csv") 17 print("[INFO] -- Data subjects' information is imported.") 18 19 return dss </pre>	<pre> 41 def creat_time_series(dt_list, act_labels, trial_codes, mode="mag", labeled=True): 42 """ 43 num_data_cols = len(dt_list) if mode == "mag" else len(dt_list)*3 44 45 if labeled: 46 dataset = np.zeros((0, num_data_cols + 7)) # 77 --> [act, code, weight, height, age, gender, trial] 47 else: 48 dataset = np.zeros((0, num_data_cols)) 49 50 ds_list = get_ds_infos() 51 52 print("[INFO] -- Creating Time-Series") 53 for sub_id in ds_list["code"]: 54 for trial in trial_codes[act_labels]: 55 frame = "A_DeviceMotion_data/" + act + "_" + str(trial) + "/" + sub_id + ".csv" 56 raw_data = pd.read_csv(frame) 57 raw_data = raw_data.drop(["Unnamed: 0"], axis=1) 58 vals = np.zeros((len(raw_data), num_data_cols)) 59 if mode == "mag": 60 vals[:, x_id] = (raw_data[axes] ** 2).sum(axis=1) ** 0.5 61 else: 62 vals[:, x_id * 3:(x_id + 1) * 3] = raw_data[axes].values 63 vals = vals[:, :num_data_cols] 64 if labeled: 65 lbls = np.array([act_id, 66 sub_id - 1, 67 ds_list["weight"][sub_id - 1], 68 ds_list["height"][sub_id - 1], 69 ds_list["age"][sub_id - 1], 70 ds_list["gender"][sub_id - 1], 71 trial 72]) + len(raw_data) 73 vals = np.concatenate((vals, lbls), axis=1) 74 dataset = np.append(dataset, vals, axis=0) 75 cols = [] 76 for axes in dt_list: 77 if mode == "raw": 78 cols += axes 79 else: 80 cols += [str(axes[0]:(-2))] 81 if labeled: 82 cols += ["act", "id", "weight", "height", "age", "gender", "trial"] 83 84 dataset = pd.DataFrame(data=dataset, columns=cols) 85 86 # id가 22인 데이터를 제거하고 별도로 저장 87 dataset_ver = dataset[(dataset["id"] == 22) (dataset["id"] == 23)] 88 dataset = dataset[(dataset["id"] != 22) & (dataset["id"] != 23)] 89 90 return dataset, dataset_ver </pre>
set_data_types 함수	
<pre> 21 def set_data_types(data_types=["userAcceleration"]): 22 """ 31 32 dt_list = [] 33 for t in data_types: 34 if t != "attitude": 35 dt_list.append([t+".x", t+".y", t+".z"]) 36 else: 37 dt_list.append([t+".roll", t+".pitch", t+".yaw"]) 38 39 return dt_list </pre>	
dataset load 과정	
<pre> act_labels = ACT_LABELS [3:4] print("[INFO] -- dataset_jog --") print("[INFO] -- Selected activities: " + str(act_labels)) trial_codes = [TRIAL_CODES[act] for act in act_labels] dt_list = set_data_types(act) dataset_jog, dataset_ver_jog = creat_time_series(dt_list, act_labels, trial_codes, mode="raw", labeled=True) print("[INFO] -- Shape of time-Series dataset:" + str(dataset_jog.shape)) print("\n") </pre>	<pre> act_labels = ACT_LABELS [2:3] print("[INFO] -- dataset_wlk --") print("[INFO] -- Selected activities: " + str(act_labels)) trial_codes = [TRIAL_CODES[act] for act in act_labels] dt_list = set_data_types(act) dataset_wlk, dataset_ver_wlk = creat_time_series(dt_list, act_labels, trial_codes, mode="raw", labeled=True) print("[INFO] -- Shape of time-Series dataset:" + str(dataset_wlk.shape)) print("\n") </pre>

2) id가 15.0, 22.0인 데이터

3) <https://github.com/mmalekzadeh/motion-sense>

creat_time_series 함수에 센서 데이터 타입과 가져오고 싶은 동작의 종류를 인자로 전달하여 호출하면 data_subjects_info.csv 파일에 저장된 사용자의 이름 정보를 읽어와 각 사용자에게 대한 센서 정보를 가져옴과 동시에 사용자의 정보를 추가하여 dataset을 구성한다. creat_time_series 함수를 사용하여 모든 데이터를 동작별로 저장하였으며 변수는 dataset_xxx 형식으로 일관성 있게 생성하였다. 이때, 마지막 두 명에 대한 데이터를 분리하여 dataset_ver_xxx 형식으로 저장하도록 하였다.

#데이터 확인

정상적으로 수행되었는지 확인해보기 위해 각 데이터셋을 출력해보았다.

dataset_wlk							dataset_ver_wlk						
act	id	weight	height	age	gender	trial	act	id	weight	height	age	gender	trial
0.0	0.0	102.0	188.0	46.0	1.0	7.0	0.0	15.0	96.0	172.0	29.0	0.0	7.0
0.0	0.0	102.0	188.0	46.0	1.0	7.0	0.0	15.0	96.0	172.0	29.0	0.0	7.0
0.0	0.0	102.0	188.0	46.0	1.0	7.0	0.0	15.0	96.0	172.0	29.0	0.0	7.0
0.0	0.0	102.0	188.0	46.0	1.0	7.0	0.0	15.0	96.0	172.0	29.0	0.0	7.0
0.0	0.0	102.0	188.0	46.0	1.0	7.0	0.0	15.0	96.0	172.0	29.0	0.0	7.0
...
0.0	23.0	74.0	173.0	18.0	0.0	15.0	0.0	22.0	68.0	170.0	25.0	0.0	15.0
0.0	23.0	74.0	173.0	18.0	0.0	15.0	0.0	22.0	68.0	170.0	25.0	0.0	15.0
0.0	23.0	74.0	173.0	18.0	0.0	15.0	0.0	22.0	68.0	170.0	25.0	0.0	15.0
0.0	23.0	74.0	173.0	18.0	0.0	15.0	0.0	22.0	68.0	170.0	25.0	0.0	15.0
0.0	23.0	74.0	173.0	18.0	0.0	15.0	0.0	22.0	68.0	170.0	25.0	0.0	15.0

dataset_wlk과 달리 dataset_ver_wlk의 경우 분리하였던 15.0, 22.0에 해당하는 id가 담겨있는 것을 확인할 수 있다.

#데이터 전처리 및 결측값 확인 함수 선언

split_dataset_by_id 함수	print_lengths_by_id_and_trial 함수
<pre> 1 def split_dataset_by_id(dataset, dataset_name): 2 # 각 id별로 데이터를 그룹화하여 pickle로 저장하는 함수 3 datasets_by_id = {} 4 5 for id, group in dataset.groupby('id'): 6 group_reset = group.reset_index(drop=True) # 인덱스를 초기화 7 datasets_by_id[id] = group_reset 8 9 return datasets_by_id </pre>	<pre> 63 def print_lengths_by_id_and_trial(datasets_by_id_and_trial): 64 # 각 id와 trial 별로 데이터프레임의 길이를 출력하는 함수 65 ids = [] 66 trials = [] 67 lengths = [] 68 69 for name, df in datasets_by_id_and_trial.items(): 70 parts = name.split('_') 71 id_part = parts[-3] 72 trial_part = parts[-1] 73 ids.append(id_part) 74 trials.append(trial_part) 75 lengths.append(len(df)) 76 77 # id, trial, 길이 정보를 문자열로 만들어 출력 78 id_str_list = [f'{id}' for id in ids] 79 trial_str_list = [f'{trial}' for trial in trials] 80 length_str_list = [str(length) for length in lengths] 81 82 id_str = " ".join(s.ljust(12) for s in id_str_list) 83 trial_str = " ".join(s.ljust(12) for s in trial_str_list) 84 length_str = " ".join(s.ljust(12) for s in length_str_list) 85 86 print(id_str) 87 print(trial_str) 88 print(length_str) 89 90 # 가장 긴 데이터프레임과 짧은 데이터프레임을 찾아 출력 91 max_length = max(lengths) 92 min_length = min(lengths) 93 max_index = lengths.index(max_length) 94 min_index = lengths.index(min_length) 95 max_id = ids[max_index] 96 max_trial = trials[max_index] 97 min_id = ids[min_index] 98 min_trial = trials[min_index] 99 100 print(f"\nLongest id: {max_id}, trial: {max_trial} with length {max_length}") 101 print(f"\nShortest id: {min_id}, trial: {min_trial} with length {min_length}") 102 </pre>
split_by_trial 함수	
<pre> 11 def split_by_trial(datasets_by_id): 12 # 각 id별로 그룹화된 데이터에서 trial 별로 다시 그룹화하여 pickle로 저장하는 함수 13 datasets_by_id_and_trial = {} 14 15 for name, data in datasets_by_id.items(): 16 id_part = name.split('_')[-1] 17 for trial, group in data.groupby('trial'): 18 key = f'{name}_{trial}_{trial}' 19 datasets_by_id_and_trial[key] = group.reset_index(drop=True) 20 21 return datasets_by_id_and_trial 22 </pre>	
nan_check 함수	
<pre> 48 def nan_check(datasets_by_id): 49 # 데이터프레임에 결측치가 있는지 확인하는 함수 50 all_no_missing = True 51 52 for name, data in datasets_by_id.items(): 53 missing_values = data.isnull().sum() 54 if missing_values.any(): 55 all_no_missing = False 56 print(f"Missing values in {name}:") 57 print(missing_values) 58 print() 59 60 if all_no_missing: 61 print("== No missing value ==") 62 </pre>	

split_dataset_by_id 함수 : 각 id 별로 데이터를 그룹화하여 딕셔너리에 저장하는 함수
split_by_trial 함수 : 각 id별로 그룹화된 데이터에서 trial 별로 다시 그룹화하여 딕셔너리에 저장하는 함수
nan_check : id와 trial로 분리된 데이터프레임에 결측치가 있는지 확인하는 함수
print_lengths_by_id_and_trial : 각 id와 trial 별로 데이터프레임의 길이를 출력하는 함수

#전처리한 데이터 확인

위와 같이 함수를 구성하고 실제 수행한 결과는 다음과 같다.

전처리 수행 및 데이터 확인 코드

```
1 print("wlk dataset")
2 dataset_wlk_by_id = split_dataset_by_id(dataset_wlk, 'dataset_wlk')
3 dataset_wlk_by_id_and_trial = split_by_trial(dataset_wlk_by_id)
4 print_lengths_by_id_and_trial(dataset_wlk_by_id_and_trial)
5 nan_check(dataset_wlk_by_id_and_trial)
6
7 print("\n\njog dataset")
8 dataset_jog_by_id = split_dataset_by_id(dataset_jog, 'dataset_jog')
9 dataset_jog_by_id_and_trial = split_by_trial(dataset_jog_by_id)
10 print_lengths_by_id_and_trial(dataset_jog_by_id_and_trial)
11 nan_check(dataset_jog_by_id_and_trial)
12
13 print("\n\njog_ver dataset")
14 dataset_ver_wlk_by_id = split_dataset_by_id(dataset_ver_wlk, 'dataset_wlk')
15 dataset_ver_wlk_by_id_and_trial = split_by_trial(dataset_ver_wlk_by_id)
16 print_lengths_by_id_and_trial(dataset_ver_wlk_by_id_and_trial)
17 nan_check(dataset_ver_wlk_by_id_and_trial)
18
19 print("\n\njog_ver dataset")
20 dataset_ver_jog_by_id = split_dataset_by_id(dataset_ver_jog, 'dataset_jog')
21 dataset_ver_jog_by_id_and_trial = split_by_trial(dataset_ver_jog_by_id)
22 print_lengths_by_id_and_trial(dataset_ver_jog_by_id_and_trial)
23 nan_check(dataset_ver_jog_by_id_and_trial)
```

전처리 수행 및 데이터 확인 결과

```
wlk dataset
id: 0.0 | id: 0.0 | id: 0.0 | id: 1.0 | id: 1.0 | id: 1.0 | id: 2.0 | id: 2.0 | id: 2.0
trial: 7.0 | trial: 8.0 | trial: 15.0 | trial: 7.0 | trial: 8.0 | trial: 15.0 | trial: 7.0 | trial: 8.0 | trial: 15.0
5439 | 4340 | 1333 | 7681 | 5614 | 6305 | 6410 | 5194 | 3594

Longest: id: 17.0, trial: 7.0 with length 8402
Shortest: id: 0.0, trial: 15.0 with length 1333
== No missing value ==

jog dataset
id: 0.0 | id: 0.0 | id: 1.0 | id: 1.0 | id: 2.0 | id: 2.0 | id: 3.0 | id: 3.0 | id: 4.0
trial: 9.0 | trial: 16.0 | trial: 9.0 | trial: 16.0 | trial: 9.0 | trial: 16.0 | trial: 9.0 | trial: 16.0 | trial: 9.0
4861 | 1567 | 4966 | 1452 | 4753 | 1216 | 5294 | 921 | 3558

Longest: id: 20.0, trial: 9.0 with length 5535
Shortest: id: 4.0, trial: 16.0 with length 765
== No missing value ==

jog_ver dataset
id: 15.0 | id: 15.0 | id: 15.0 | id: 22.0 | id: 22.0 | id: 22.0
trial: 7.0 | trial: 8.0 | trial: 15.0 | trial: 7.0 | trial: 8.0 | trial: 15.0
7366 | 5446 | 2744 | 6505 | 5441 | 3306

Longest: id: 15.0, trial: 7.0 with length 7366
Shortest: id: 15.0, trial: 15.0 with length 2744
== No missing value ==

jog_ver dataset
id: 15.0 | id: 15.0 | id: 22.0 | id: 22.0
trial: 9.0 | trial: 16.0 | trial: 9.0 | trial: 16.0
4814 | 1699 | 4815 | 1147

Longest: id: 22.0, trial: 9.0 with length 4815
Shortest: id: 22.0, trial: 16.0 with length 1147
== No missing value ==
```

위 이미지는 print_lengths_by_id_and_trial 함수와, nan_check 함수를 사용해 출력한 각 데이터프레임 정보 중 일부이다. 현재 데이터는 어떤 동작을 취하는지에 따라 데이터의 길이가 다르며, 무엇보다 개개인의 차이에 의해 똑같은 동작을 수행하여도 그 길이가 다르다.

또한, 각 데이터에 대해 결측치가 있는지 확인해보았으나 결측치는 존재하지 않는 것을 확인할 수 있었다.

#그래프 표출 함수

plot_line 함수	plot_heatmap 함수
<pre> 1 def plot_line(datasets, data_name, id, trial, data_type, x_end=None, x_start=None): 2 # 데이터 이름 생성 3 name_to_plot = f'{data_name}_{id}_{trial}_{trial}' 4 5 # 데이터 선택 및 특정 열 선택 6 if name_to_plot in datasets: 7 columns_to_plot = [col.strip() for col in data_type.split(',')] 8 data_to_plot = datasets[name_to_plot][columns_to_plot] 9 10 11 12 # 라인 차트 출력 13 data_to_plot.plot(title=f'Dataset for ID {id} and Trial {trial} in {data_name}', figsize=(14, 4)) 14 15 # x_start와 x_end가 지정된 경우 해당 범위의 데이터 선택 16 if x_start is not None and x_end is not None: 17 plt.xlim(x_start, x_end) 18 elif x_end is not None: 19 plt.xlim(0, x_end) 20 21 plt.xlabel('Index') 22 plt.ylabel('Values') 23 plt.show() 24 25 print(f"ID {id} with Trial {trial} not found in the datasets of type {data_name}.") </pre>	<pre> 35 def plot_heatmap(datasets, data_name, id, trial, data_type): 36 # 데이터 이름 생성 37 name_to_plot = f'{data_name}_{id}_{trial}_{trial}' 38 39 # 데이터 선택 및 특정 열 선택 40 if name_to_plot in datasets: 41 columns_to_plot = [col.strip() for col in data_type.split(',')] 42 data_to_plot = datasets[name_to_plot][columns_to_plot] 43 44 # 데이터 타입 변환 및 결측값 처리 45 data_to_plot = data_to_plot.apply(pd.to_numeric, errors='coerce') 46 data_to_plot = data_to_plot.dropna() 47 48 if data_to_plot.empty: 49 print(f"No valid numeric data available for {name_to_plot}") 50 return 51 52 # 상관관계 계산 53 corr = data_to_plot.corr() 54 55 # 히트맵 출력 56 fig, ax = plt.subplots(figsize=(10, 8)) 57 im = ax.imshow(corr.values, cmap='coolwarm') 58 59 # 라벨 선택 60 ax.set_xticks(np.arange(len(corr.columns))) 61 ax.set_yticks(np.arange(len(corr.columns))) 62 ax.set_xticklabels(corr.columns) 63 ax.set_yticklabels(corr.columns) 64 65 # 논문 라벨 회전 및 정렬 설정 66 plt.setp(ax.get_xticklabels(), rotation=45, ha="right", 67 rotation_mode="anchor") 68 69 # 데이터 차원별로 텍스트 주석 생성 70 for i in range(len(corr.columns)): 71 for j in range(len(corr.columns)): 72 text = ax.text(j, i, np.round(corr.iloc[i, j], decimals=2), 73 ha="center", va="center", color="black") 74 75 # 컬러바 추가 76 cbar = ax.figure.colorbar(im, ax=ax, cmap='coolwarm') 77 cbar.ax.set_ylabel('Correlation', rotation=-90, va="bottom") 78 79 plt.title(f'Heatmap for ID {id} and Trial {trial} in {data_name}') 80 plt.show() 81 82 print(f"ID {id} with Trial {trial} not found in the datasets of type {data_name}.") </pre>
plot_boxplot 함수	
<pre> 84 def plot_boxplot(datasets, data_name, id, trial, data_type): 85 # 데이터 이름 생성 86 name_to_plot = f'{data_name}_{id}_{trial}_{trial}' 87 88 # 데이터 선택 및 특정 열 선택 89 if name_to_plot in datasets: 90 columns_to_plot = [col.strip() for col in data_type.split(',')] 91 data_to_plot = datasets[name_to_plot][columns_to_plot] 92 93 # 데이터 타입 변환 및 결측값 처리 94 data_to_plot = data_to_plot.apply(pd.to_numeric, errors='coerce') 95 data_to_plot = data_to_plot.dropna() 96 97 if data_to_plot.empty: 98 print(f"No valid numeric data available for {name_to_plot}") 99 return 100 101 # 박스플롯 출력 102 plt.figure(figsize=(8, 4)) 103 sns.boxplot(data=data_to_plot) 104 plt.title(f'Boxplot for ID {id} and Trial {trial} in {data_name}') 105 plt.xlabel('Variables') 106 plt.ylabel('Values') 107 plt.show() 108 109 print(f"ID {id} with Trial {trial} not found in the datasets of type {data_name}.") </pre>	

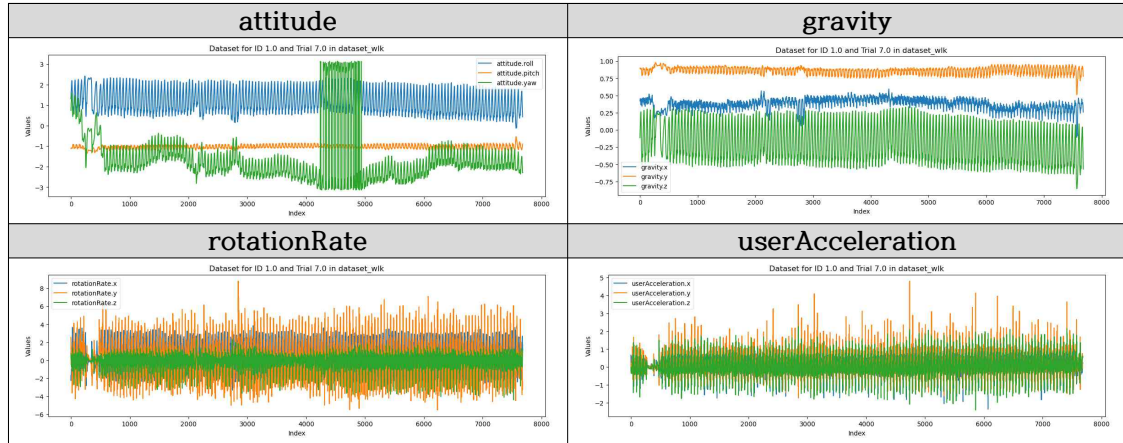
line chart, boxplot, heatmap 시각화 도구를 사용하여 데이터에 대한 분석을 진행하기 위해 각 동작을 함수로 구현하였다. 모든 함수는 출력할 데이터프레임, 데이터 이름, id, trial 그리고 출력할 센서 타입을 입력받으며, 예외적으로 plot_line 함수는 2개의 인자를 더 받는다. plot_line의 경우 x_end와 x_start 인자를 받으며 이를 통해서 시각화하고 싶은 데이터의 길이를 조정할 수 있다. 이때, 두 인자는 입력받지 않을 시 모든 데이터를 출력하며, 첫 번째 인자, 즉 x_end를 입력받으면 시작부터 x_end까지의 데이터를 출력한다.

위와 같이 함수에 대한 구성을 마쳤으므로 본격적으로 wlk 동작에 대한 분석할 것이다. 먼저, 각 데이터 타입 간 비교를 수행한다.

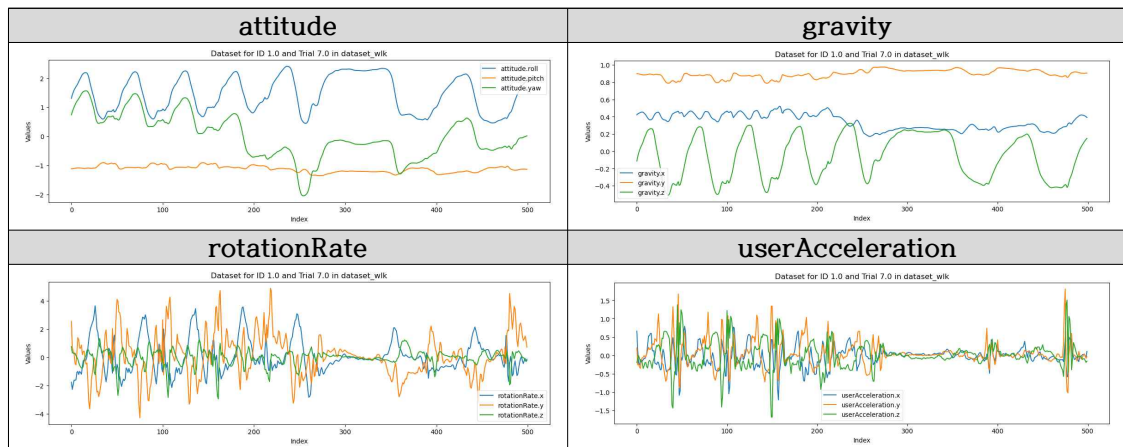
#Datatype 간 비교 (걷기)

데이터 타입 간의 비교를 위해 걷기 동작에서 id 1.0 trail 7.0의 데이터를 시각화하여 분석하였다.

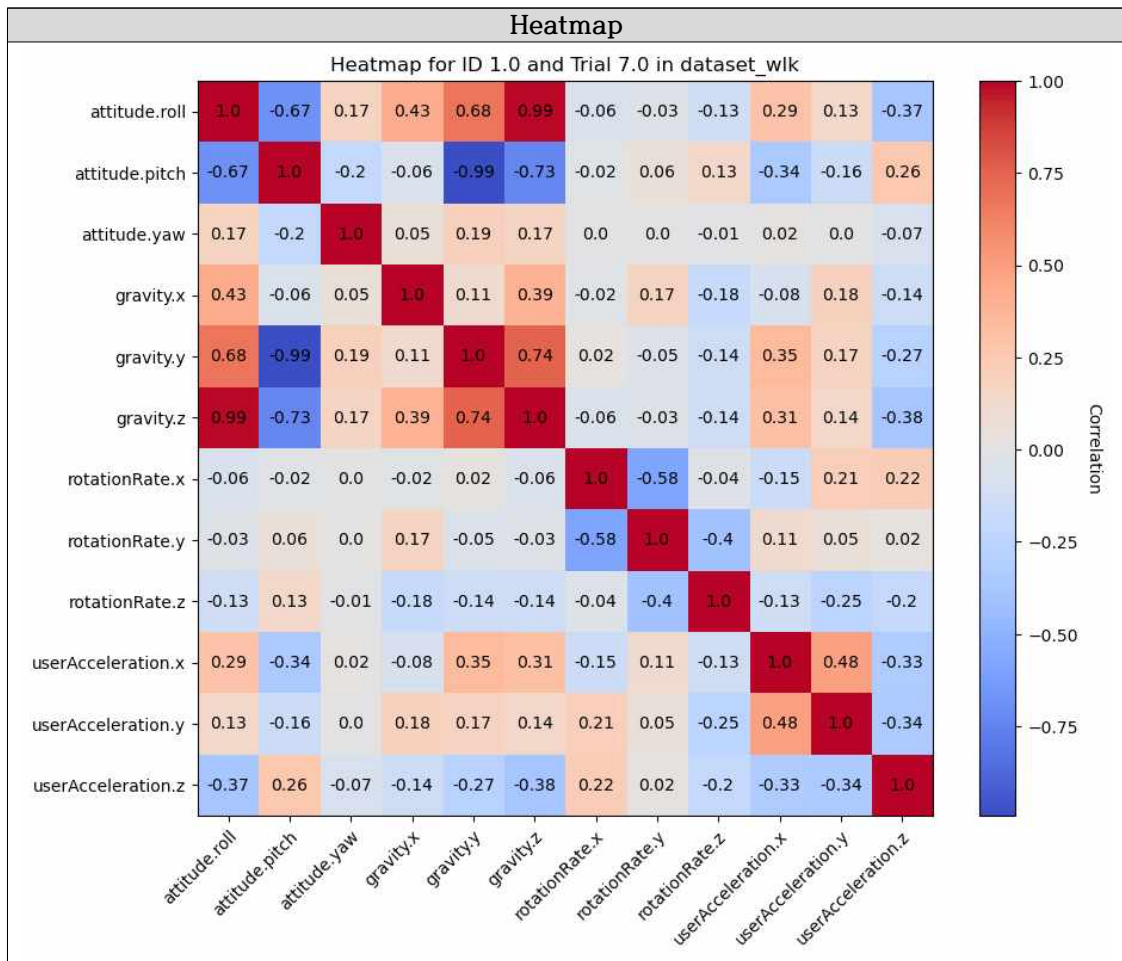
다음은 라인차트의 결과이다.



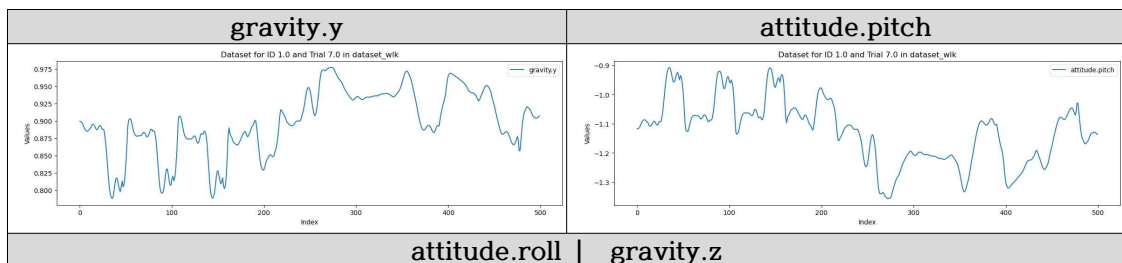
위는 모든 센서 타입별로 데이터를 출력해본 결과이다. 데이터의 길이가 약 8000으로 그래프의 밀도가 높아 구체적인 분석은 힘들다. 모든 센서 타입 데이터들은 진동하는 것으로 보이나 구체적인 분석이 힘들므로 0부터 500까지의 데이터만 출력해보았다. 아래는 그 결과이다.

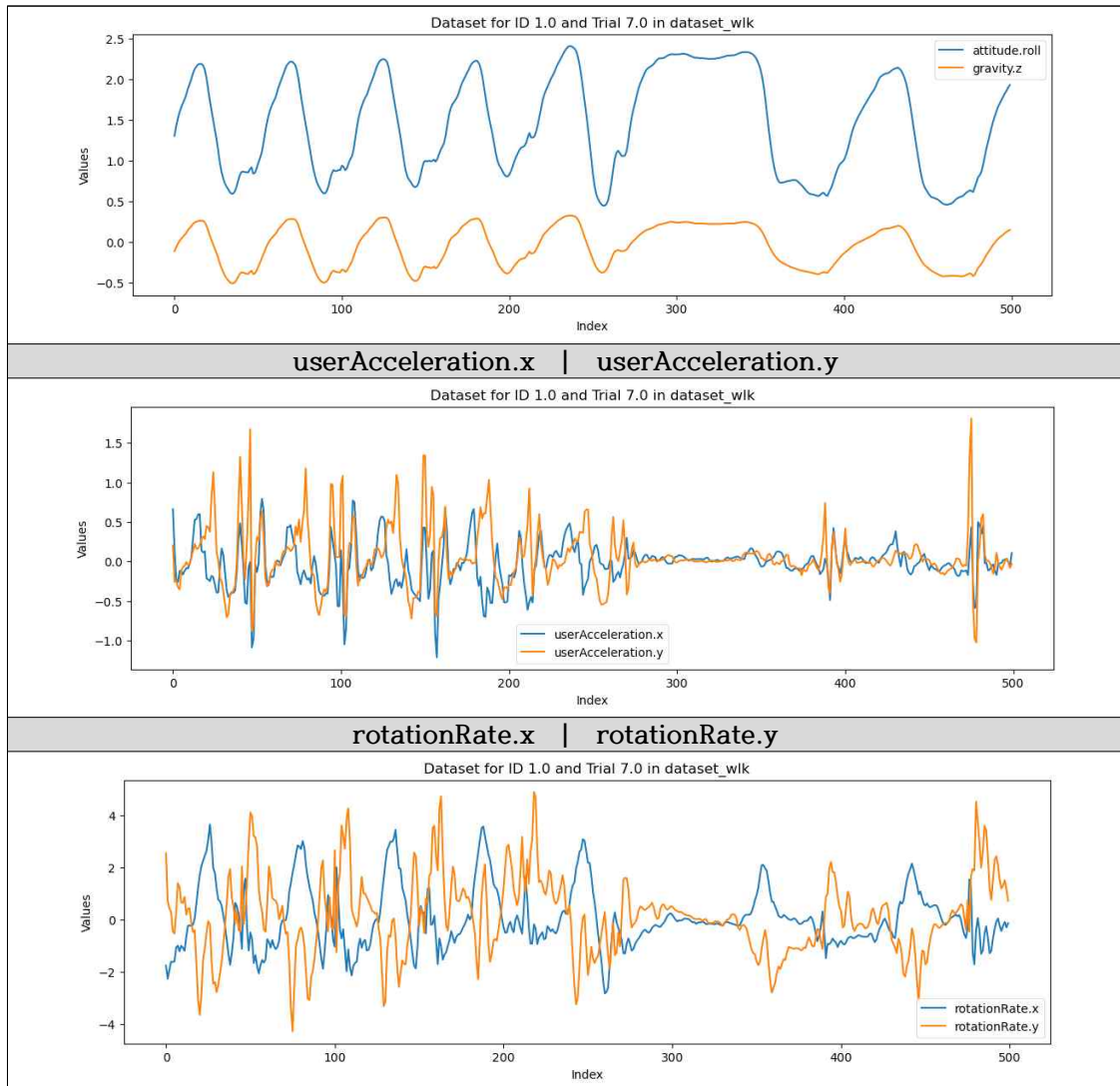


attitude와 gravity 그래프는 어느 정도 눈에 띄는 규칙성을 가진 진동 운동을 하는 것처럼 보인다. 현재 분석 중인 데이터가 걸음걸이에 대한 데이터이므로 attitude와 gravity 그래프에서 주기성이 걸음에서 비롯되었다고 유추할 수 있다. 특히 attitude 데이터의 경우 roll과 yaw 그래프가 매우 비슷한 형태를 가진다. 반면, gravity 데이터는 gravity.z에서만 주기성이 보이며 나머지 데이터는 데이터간 상관관계를 파악하기 힘들다. rotationRate와 userAcceleration 데이터도 규칙적으로 움직이는 것처럼 보이지만 노이즈가 많아 그 규칙성을 확인하기 힘들다. 각 그래프간 데이터간의 상관 관계는 다음 heatmap을 통해 구체적으로 확인해보고자 한다. 다음은 heatmap으로 모든 센서 타입간 상관관계를 확인해본 결과이다.



예상과 달리 attitude.roll과 attitude.yaw 데이터는 line chart에서 매우 비슷해 보였으나 0.17로 낮은 상관관계를 보였다. 반대로 attitude.pitch는 -0.67로 상당히 높은 상관관계를 보였다. 또한, attitude.roll과 gravity.z는 0.99로 거의 일치하는 것을 볼 수 있다. 오히려 attitude.pitch는 gravity.y와 -0.99로 마찬가지로 정반대의 상관관계를 보인다. heatmap 결과의 다른 특징은 attitude와 gravity는 어느 정도의 높은 관계성을 보여주나 rotationRate와 userAcceleration과는 전체적으로 매우 낮은 상관관계를 가진다. 또한, rotationRate와 userAcceleration은 둘 사이에도 낮은 상관관계를 가진다. 다만, 각자의 x, y, z 사이에서는 약간의 상관관계가 있는 것을 확인할 수 있다. 이를 눈으로 확인해보기 위해 그래프로 어느 정도의 상관관계를 가진 데이터를 하나의 그래프에 출력하여 확인해보았다.

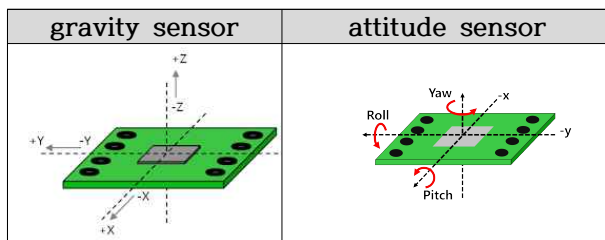




결과 gravity.y와 attitude.pitch가 확실해 반대의 상관관계를 가진 것을 확인할 수 있었으며 attitude.roll과 gravity.z 또한 높은 상관관계를 가진 것을 볼 수 있다.

userAcceleration.x와 userAcceleration.y 그리고 rotationRate.x와 rotationRate.y 또한 이전 모두 line chart로 표현했을 때 보다 유의미한 상관관계가 있는 것을 확인 가능했다.

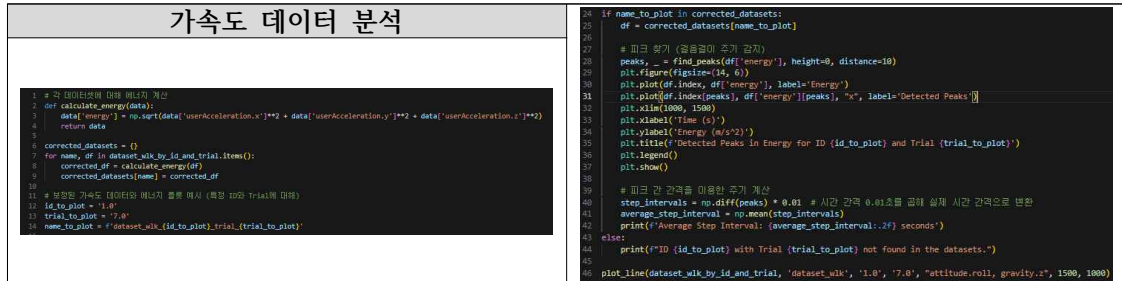
gravity.y와 attitude.pitch 그리고 attitude.roll과 gravity.z가 높은 상관관계를 보이는 이유는 다음과 같이 예측할 수 있다.



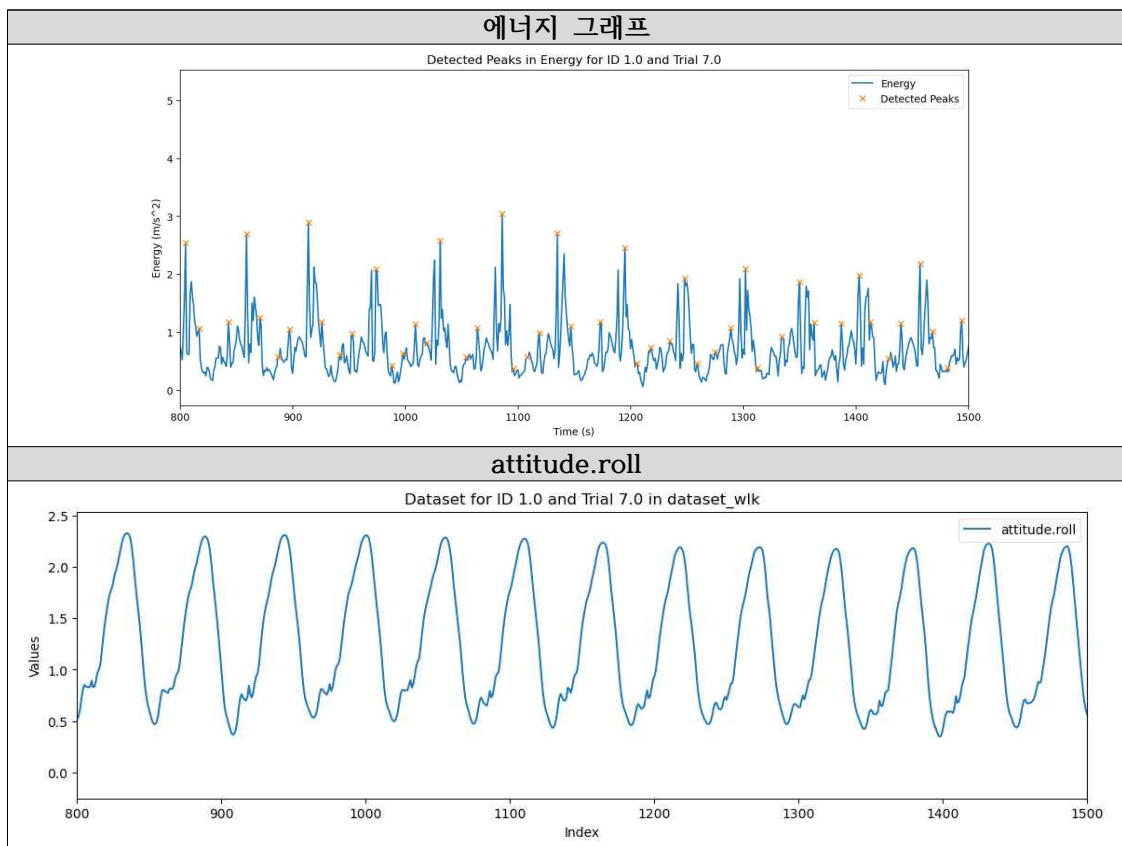
좌측 이미지에서 pitch와 roll의 값을 증가시켜 핸드폰의 기울기를 변경할 때 gravity sensor의 y와 x의 값이 각각 pitch와 roll의 움직임에 비례 혹은 반비례할 것이다.

따라서, gravity sensor나 attitude sensor를 택할 경우 attitude.yaw와 gravity.z 값 중 걸음걸이 움직임의 특성을 더욱 잘 담은 센서 타입을 선택하는 것이 좋을 것이다.

일반적으로 핸드폰에 있는 만보기 기능은 가속도계 데이터를 사용하므로 가속도계 데이터를 사용하는 것으로 알려져 있다. 가속도 센서 데이터를 사용하여 걸음걸이를 유추하기 위해 가속도 센서 데이터를 에너지 데이터로 변경하여 데이터를 분석해 보았다.

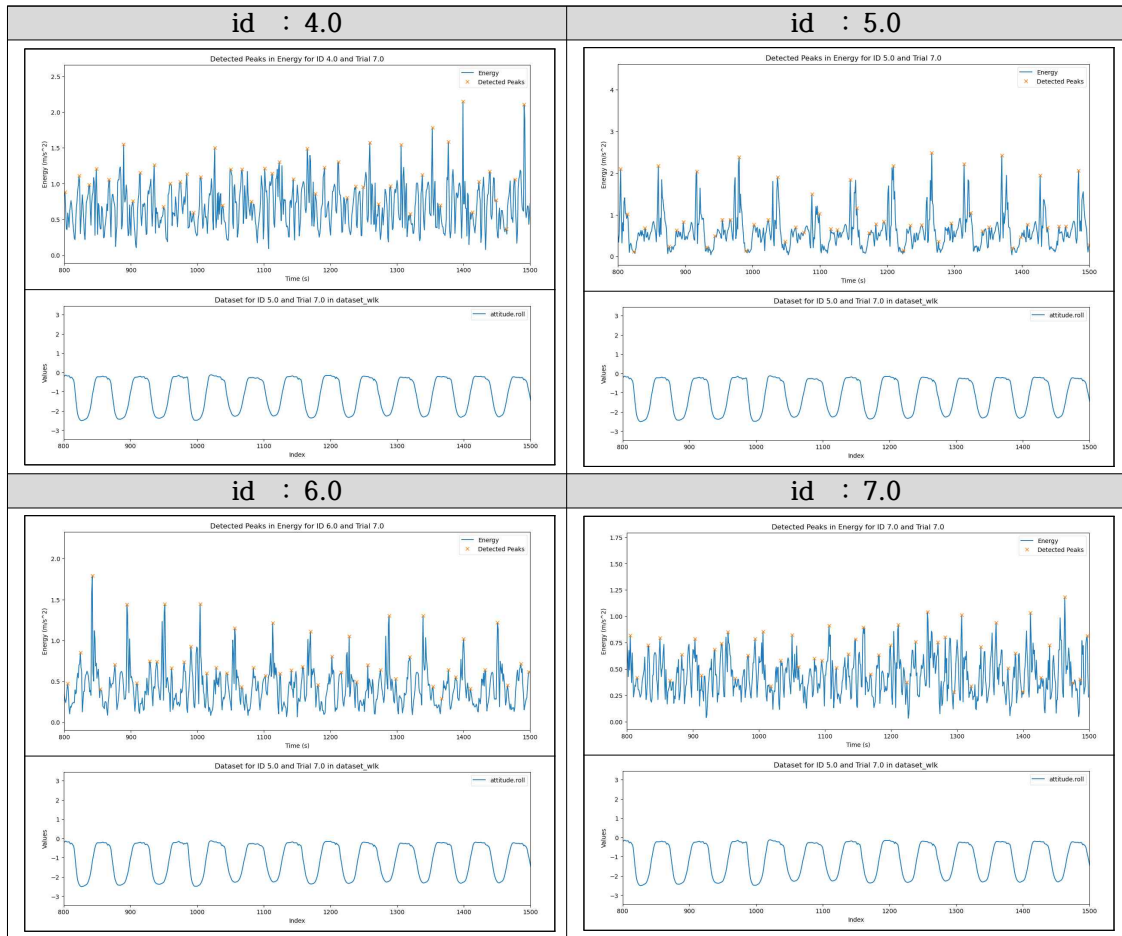


제공된 문서에 따르면 userAcceleration 데이터는 이미 중력 성분을 제거한 데이터이다. 따라서 $E = \sqrt{x^2 + y^2 + z^2}$ 를 통하여 에너지를 계산하고 그래프로 나타내 보았다.



변환된 에너지 그래프는 여전히 노이즈가 많다. 이 때문에 많은 피크가 감지된 것을 확인할 수 있다. 필터링을 통해 개선의 여지가 있지만 주목할 점은 에너지 그래프의 피크와 attitude.roll의 피크가 매우 유사한 모습을 보인다. 그렇다면 복잡한 과정 없이 attitude.roll 데이터를 사용해서 걸음걸이를 파악해도 충분히 신뢰할 만한 주기를 얻을 수 있을 것이다.

그렇다면 항상 이러한 결과가 나타는지 다른 사람에 대해서도 확인을 해 보았다.

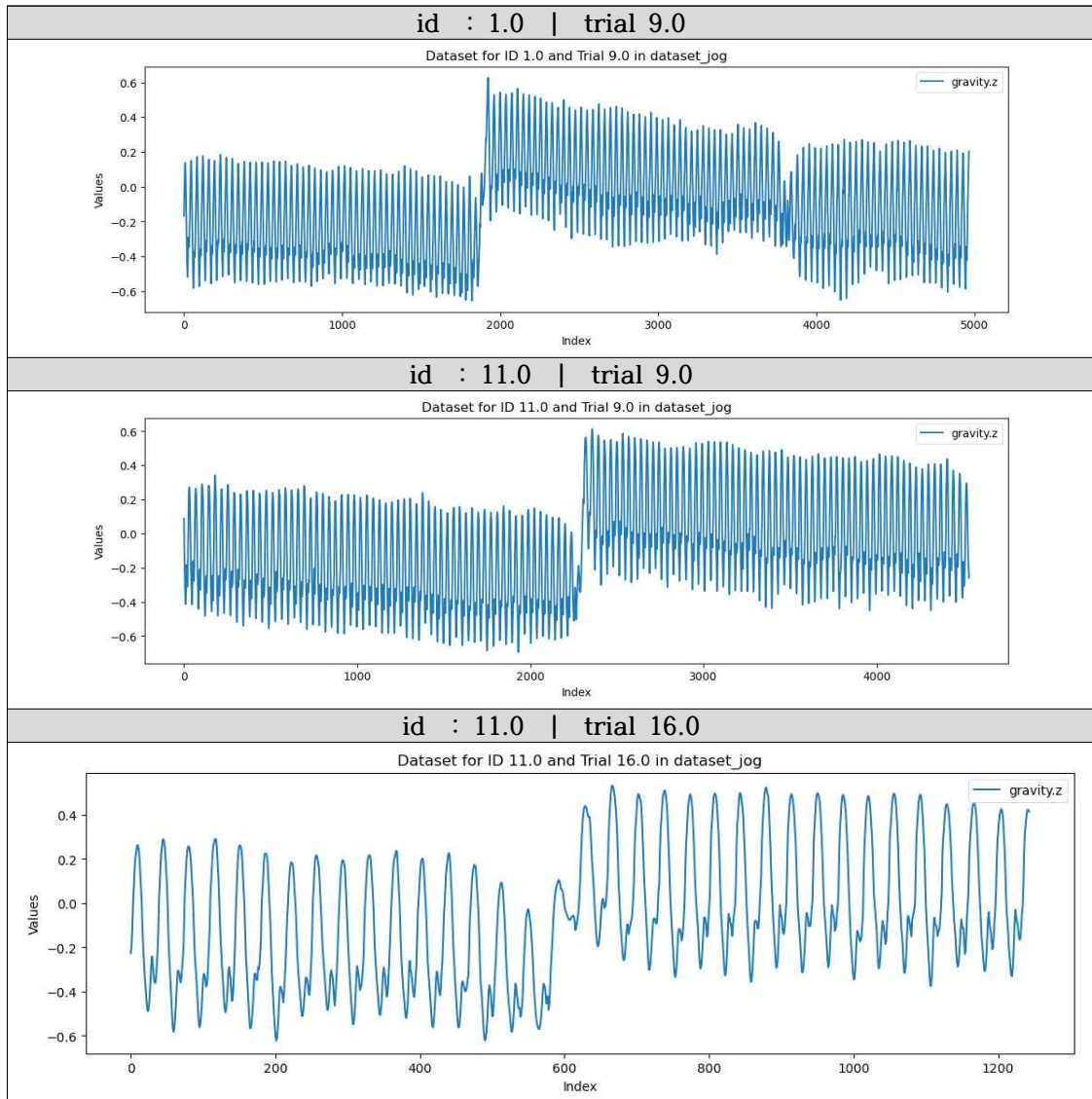


id 5와 6은 매우 뚜렷하게 에너지 그래프와 attitude.roll 그래프와 유사한 것을 볼 수 있다. 하지만, 4와 7은 노이즈가 커 뚜렷하게 보이진 않는다. 하지만, 전반적으로 attitude.roll과 유사한 것을 볼 수 있다.

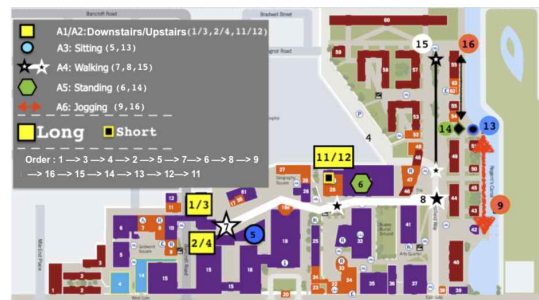
따라서, attitude.roll과 매우 높은 유사도를 가지는 gravity.z, 혹은 attitude.pitch나 gravity.y를 사용하여 걸음걸이 데이터를 추출하고 신장별 특징을 확인할 것이다.

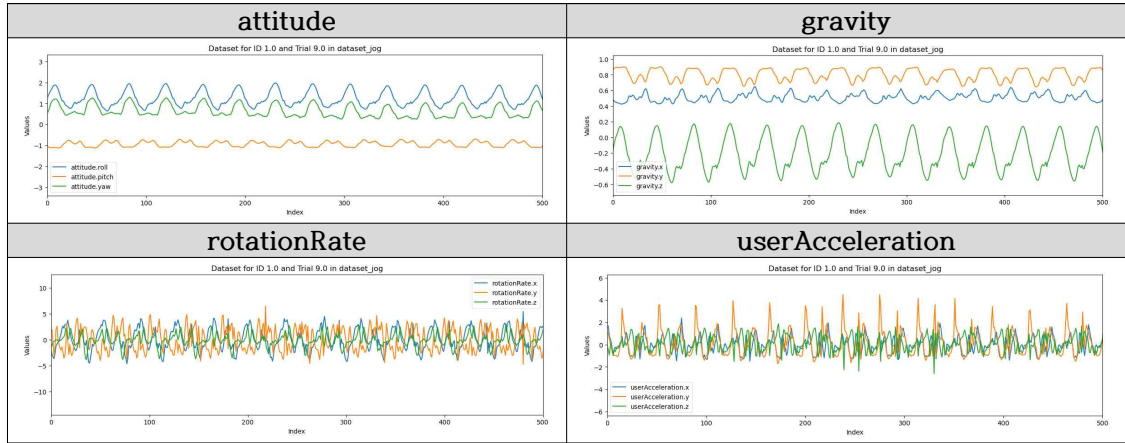
#Datatype 간 비교 (조깅)

걷기 과정에서 분석한 결과를 토대로 조깅에 대한 정보도 시각화하여 분석해보았다.

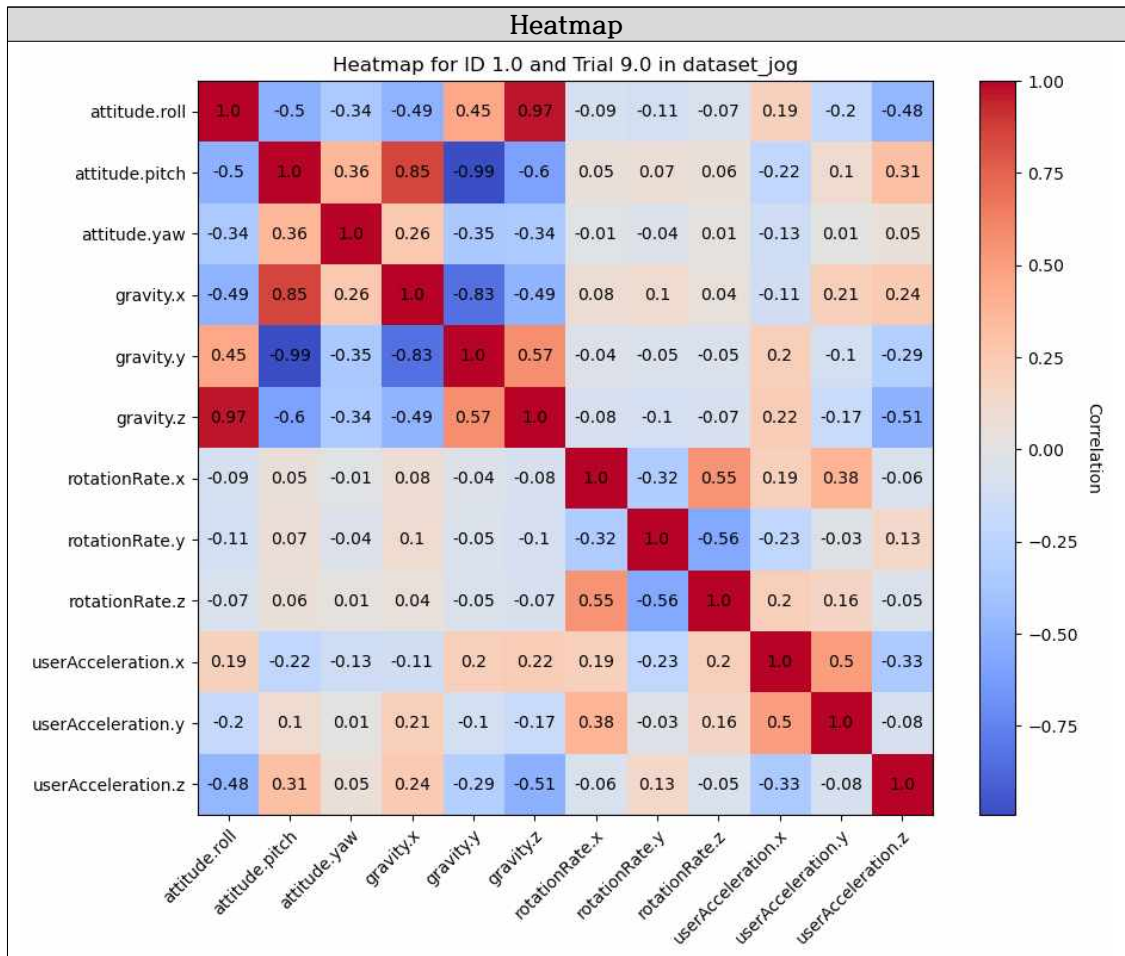


조깅 데이터의 경우 한번 출력되는 경우가 존재한다. 이유를 찾기 위해 공식 문서를 찾아본 결과 조깅 동작은 왕복 운동을 수행한 것으로 확인할 수 있었다. 조깅에 대한 두 가지 trial 9.0과 16.0 모두 왕복 동작으로 데이터가 출력되는 모습을 관측할 수 있었다. 이후 데이터를 확대하여 구체적으로 확인해보았다.



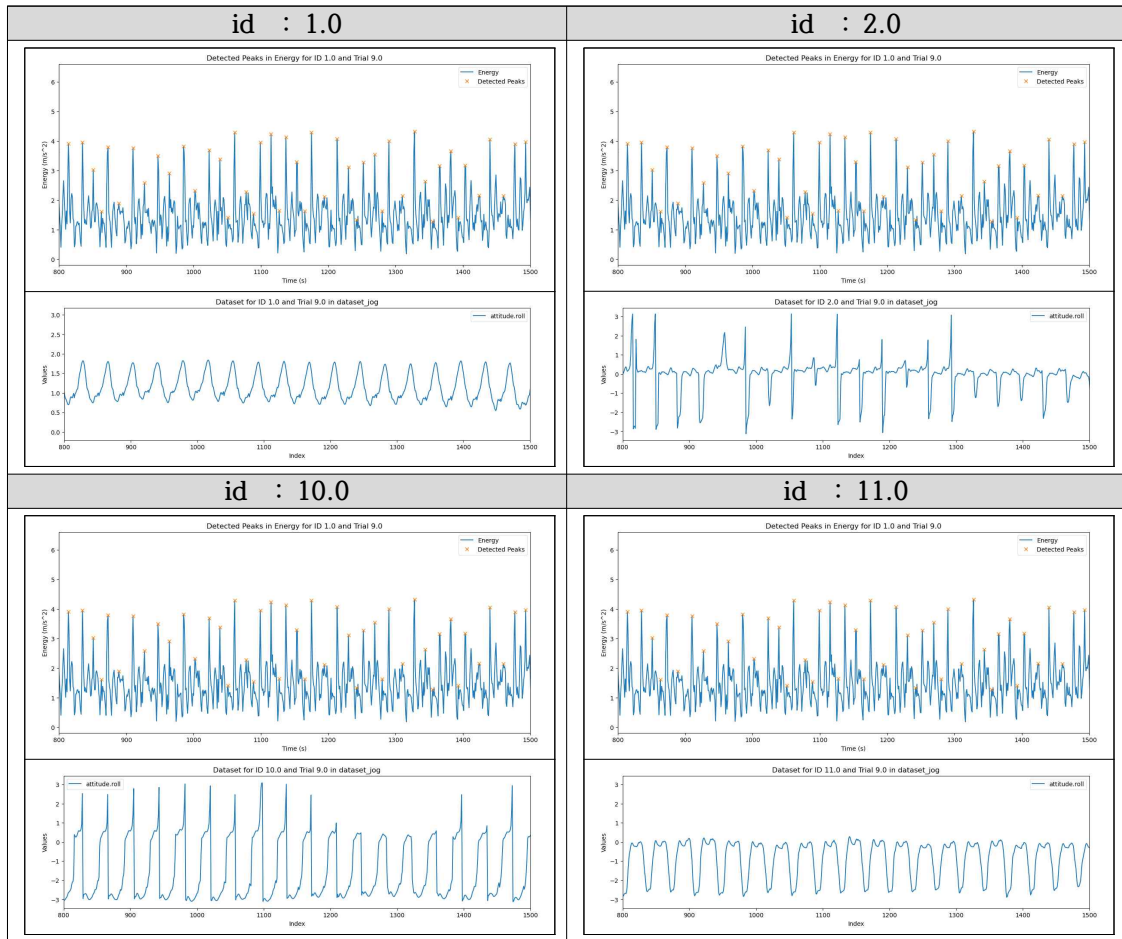


id 1의 trial 9 조깅 운동에 대해 라인차트를 확인해보니 이전 walk와 크게 차이를 보이지 않았다. 따라서 바로 heatmap을 확인해 보았다.

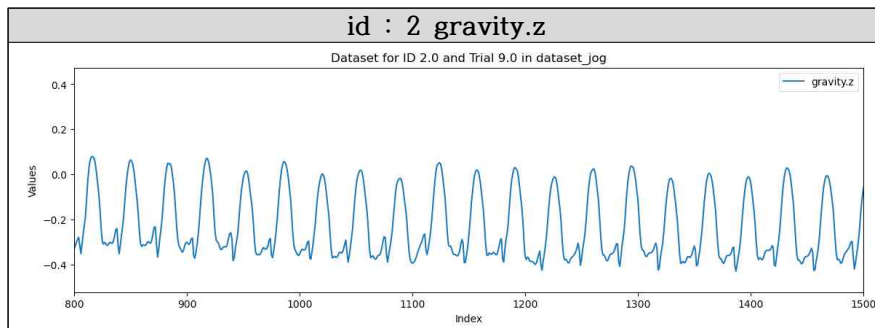


heatmap의 경우도 마찬가지로 이전과 매우 유사하다. 주목할 만한 점은 이전과 마찬가지로 gravity.y와 attitude.pitch 그리고 attitude.roll과 gravity.z가 매우 높은 상관관계를 보이는 점, 그리고 이전과 달리 attitude.pitch와 gravity.x가 -0.06에서 0.85로 상관관계 수치가 매우 높아졌다는 점이다. 하지만, 이전 걷기 동작에서 선택했던 attitude.roll, attitude.pitch,

gravity.y와 gravity.x는 여전히 유효하다. 마지막으로 가속도 센서 데이터를 에너지로 변환하여 확인해본 결과이다.



전반적으로 에너지 그래프와 attitude.roll은 일치하였으나 id 2에서 다소 문제가 될 수 있는 부분을 발견하였다. 데이터에 피크를 일정하게 추출할 수 있어야 걸음걸이를 확인할 수 있으나 이상한 노이즈가 포함되어 있는 것을 확인할 수 있었다. 따라서, 다른 후보군이던 gravity.z를 사용하여 확인해본 결과 다음과 같다.



다행이도 gravity.z의 경우 정상적인 피크를 가지고 있는 것을 확인할 수 있었다. 따라서, gravity.z 데이터를 사용하여 1차적으로 신장별 걸음걸이를 분석할 것이다.

#데이터 가공 및 주기 확인

데이터를 가공하는 과정은 다음과 같다.

1. gravity.z 데이터에서 피크값을 찾고 피크 사이의 간격 즉 주기를 추출한다.
2. gravity.z의 주기를 바탕으로 energy 그래프를 적분한다.
3. 주기별 적분된 데이터의 평균을 내고 평균에 가까운 80%의 데이터를 추출한다.
4. 추출한 80%의 데이터에 대해서 다시 평균을 추출

위와 같은 과정을 통해 가속도 데이터에 대한 필터링 없이 쉽게 에너지에 대한 정보를 추출하고, 추출한 데이터 중 이상치를 제거하여 개개인에 대한 평균적인 걸음걸이를 얻는다.

위 과정을 위해 선언한 함수들은 다음과 같다.

<h3>plot_line_peak</h3> <pre>1 def plot_line_peak(datasets, data_name, id, trial, data_type, x_end=None, x_start=None): 2 # 데이터 이름 설정 3 name_to_plot = f'{data_name}_{id}_trial_{trial}' 4 5 # 데이터 선택 및 특정 열 선택 6 if name_to_plot in datasets: 7 columns_to_plot = [col.strip() for col in data_type.split(',')] 8 data_to_plot = datasets[name_to_plot][columns_to_plot] 9 10 # 피크 찾기 11 peaks, _ = find_peaks(data_to_plot.iloc[:, 0], distance=10, prominence=0.15, width=1) 12 13 # 라인 차트 출력 14 plt.figure(figsize=(14, 4)) 15 plt.plot(data_to_plot.index, data_to_plot.iloc[:, 0], label='Data') 16 plt.plot(data_to_plot.index[peaks], data_to_plot.iloc[peaks, 0], "x", label='Peaks') 17 18 # x_start와 x_end가 지정한 경우 해당 범위의 데이터 선택 19 if x_start is not None and x_end is not None: 20 plt.xlim(x_start, x_end) 21 elif x_end is not None: 22 plt.xlim(0, x_end) 23 24 plt.title(f'Dataset for ID {id} and Trial {trial} in {data_name}') 25 plt.xlabel('Index') 26 plt.ylabel('Values') 27 plt.legend() 28 plt.show() 29 else: 30 print(f'ID {id} with Trial {trial} not found in the datasets of type {data_name}.')</pre>	<h3>integrate_energy_over_intervals</h3> <pre>35 # 피크 간 적분된 에너지 계산 함수 (gravity.z 사용) 36 def integrate_energy_over_intervals(data, peaks): 37 integrated_values = [] 38 for start, end in zip(peaks[:-1], peaks[1:]): 39 integrated_value = np.trapz(data['energy'], (start:end)) 40 integrated_values.append(integrated_value) 41 return np.array(integrated_values)</pre>
<h3>process_and_select_data</h3> <pre>44 # 데이터 처리 및 데이터 선택 45 def process_and_select_data(datasets): 46 selected_data = {} 47 for name, data in datasets.items(): 48 selected_values = select_closest_to_mean(data) 49 selected_data[name] = np.mean(selected_values) 50 return selected_data</pre>	<h3>select_closest_to_mean</h3> <pre>51 # 데이터 처리, 가까운 값을 선택하는 함수 52 def select_closest_to_mean(data, percentage=0.8): 53 mean_value = np.mean(data) 54 distances = np.abs(data - mean_value) 55 sorted_indices = np.argsort(distances) 56 cutoff = int(len(data) * percentage) 57 selected_indices = sorted_indices[:cutoff] 58 return data[selected_indices]</pre>
<h3>calculate_peak_intervals</h3> <pre>59 # 피크 위치 및 주기 계산 60 def calculate_peak_intervals(data): 61 peak_intervals = {} 62 for name, df in data.items(): 63 peaks, _ = find_peaks(df['gravity.z'], distance=10, prominence=0.15, width=1) 64 if len(peaks) > 1: 65 intervals = np.diff(peaks) 66 intervals = select_closest_to_mean(intervals, percentage=0.8) 67 mean_interval = np.mean(intervals) 68 peak_intervals[name] = mean_interval 69 else: 70 peak_intervals[name] = None 71 return peak_intervals</pre>	<h3>process_datasets</h3> <pre>61 # 데이터 처리 함수 62 def process_datasets(datasets): 63 processed_data = {} 64 for name, df in datasets.items(): 65 df = calculate_energy(df) 66 peaks, _ = find_peaks(df['gravity.z'], distance=10, prominence=0.15, width=1) 67 integrated_values = integrate_energy_over_intervals(df, peaks) 68 processed_data[name] = integrated_values 69 return processed_data</pre>
	<h3>calculate_energy</h3> <pre>72 # 에너지 계산 함수 73 def calculate_energy(data): 74 data['energy'] = np.sqrt(data['userAcceleration.x']**2 + data['userAcceleration.y']**2 + data['userAcceleration.z']**2) 75 return data</pre>

plot_line_peak : 이전 라인차트 함수에서 피크를 찍는 것을 추가한 함수

calculate_energy : 가속도 데이터로 에너지를 계산하는 함수

integrate_energy_over_intervals : 피크 위치를 토대로 에너지를 적분하는 함수

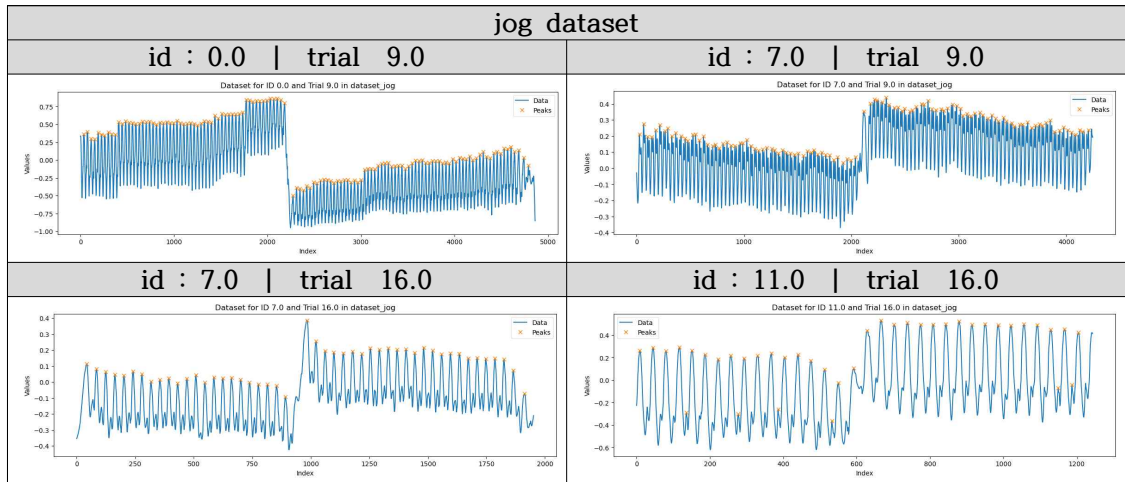
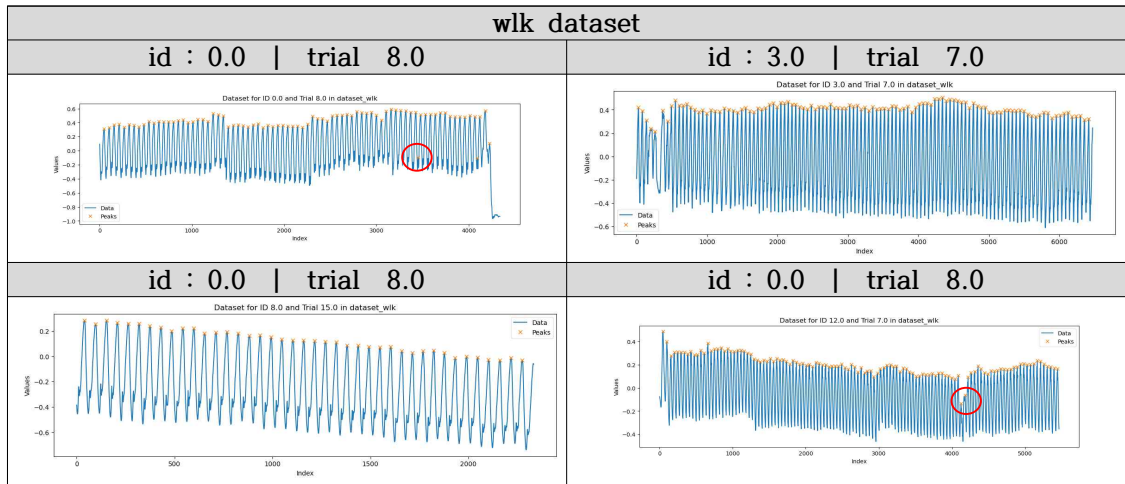
process_datasets : gravity.z의 피크를 찾고 적분을 하며, 해당 데이터를 반환하는 함수

select_closest_to_mean : 적분 값 중 평균에 가까운 80%의 데이터만 남기는 함수

process_and_select_data : 남은 80%의 데이터에서 다시 평균을 구하여 반환하는 함수

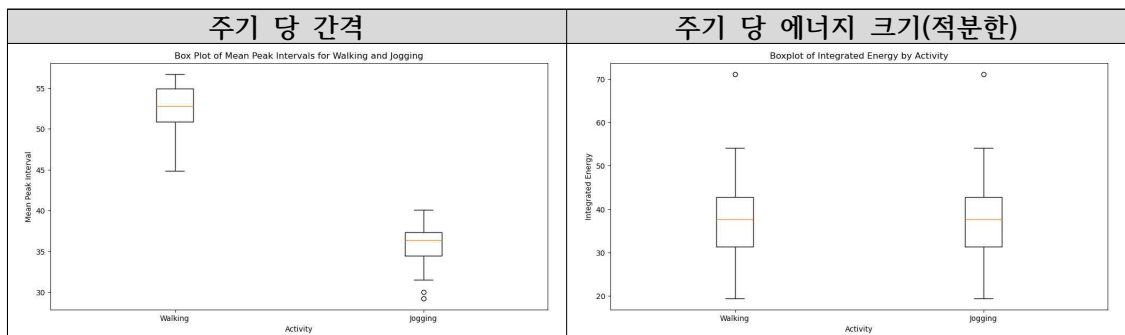
calculate_peak_intervals : 피크의 간격을 추출하여 평균에 80%에 해당하는 평균을 반환하는 함수

gravity.z 가 정말 주기를 잘 담고 있는지 확인하기 위해 모든 gravity.z 그래프를 출력하여 확인해보았다. 다음은 해당 결과 중 일부이다.

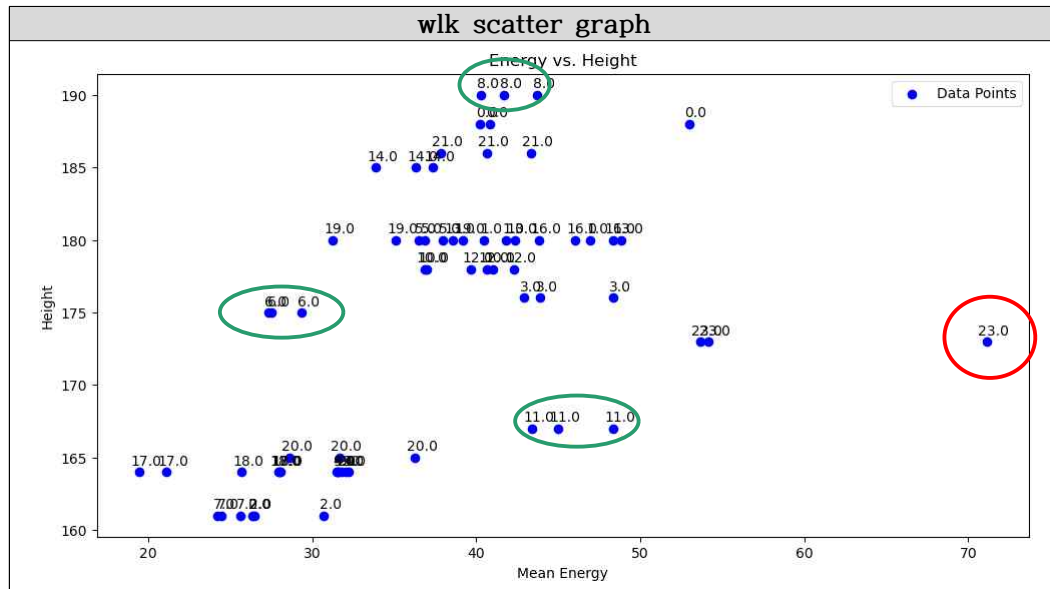


피크를 찾을 때, “distance=10, prominence=0.15, width=1”와 같은 옵션을 주었기 때문에 피크를 잘 찾았으나 걷기 데이터에서 일반적이지 않은 위치에서 피크가 존재하는 것을 확인할 수 있다. 하지만, 이러한 데이터는 평균을 내고 80%의 데이터만 택하는 과정에서 제거할 수 있을 것이다. 또한, 조깅 데이터에서 반환점에서 데이터가 출렁일 때 발생하는 오차 또한 제거할 수 있다.

다음은 동작 별 주기 간격의 평균과 적분한 에너지 크기의 평균을 box plot으로 시각화한 결과이다.



위 데이터는 가속도 에너지 데이터를 적분하는 것만으로는 걷기와 뛰기를 구분할 수 없다는 것을 의미한다. 반면, 주기의 간격 데이터는 걷기와 조깅이 겹치지 않고 잘 분리되어 있는 것을 확인할 수 있다.



위 그래프는 적분하여 얻은 결과의 평균을 scatter 그래프로 나타낸 결과이다. 위 결과를 보면 신장이 커질수록 에너지가 커지는 듯한 움직임을 보이긴 하지만, 일부 데이터는 선형적인 움직임을 따르지 않는다. 빨간 원으로 표시한 데이터는 모델을 학습하는 데 있어서 큰 방해가 될 것으로 예상된다. 하지만, 이러한 데이터는 모델 학습 단계에서 제거하여 높은 정확도를 가진 모델을 얻을 수 있을 것이다.

반면, 매우 이상적으로 측정되었음을 알 수 있는 데이터도 존재한다. 초록색 원으로 표시된 데이터들은 같은 id, 즉, 같은 사람으로부터 수집된 데이터가 서로 근접해 있는 것을 확인할 수 있다. 이는, 신장과 에너지가 비례하는 움직임을 떠나서 수집된 데이터가 일관성 있게 측정되었다는 것을 알 수 있다.

#모델 학습(wlk)

본격적으로 모델을 학습하기에 앞서 걷기 데이터에 대해서만 학습해보았다.

```
def polynomial_regression_analysis(X, y, ids, degree=1, test_size=0.2, num_iterations=1000):
    # 데이터 스케일링
    scaler = StandardScaler()
    X_scaled = scaler.fit_transform(X.reshape(-1, 1))

    r2_scores = []
    mse_scores = []
    models = []

    for _ in range(num_iterations):
        # 데이터 분할
        X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=test_size, random_state=np.random.randint(10000))

        # 다항 회귀 모델 학습
        polynomial_features = PolynomialFeatures(degree=degree)
        model = make_pipeline(polynomial_features, LinearRegression())
        model.fit(X_train, y_train)
        models.append(model)

        # 예측 및 평가
        y_pred = model.predict(X_test)
        mse = mean_squared_error(y_test, y_pred)
        r2 = r2_score(y_test, y_pred)
        mse_scores.append(mse)
        r2_scores.append(r2)

    # 결과 출력
    best_index = np.argmax(r2_scores)
    best_model = models[best_index]

    print(f"Best R^2 Score: {r2_scores[best_index]:.2f}")
    print(f"Corresponding Mean Squared Error: {mse_scores[best_index]:.2f}")

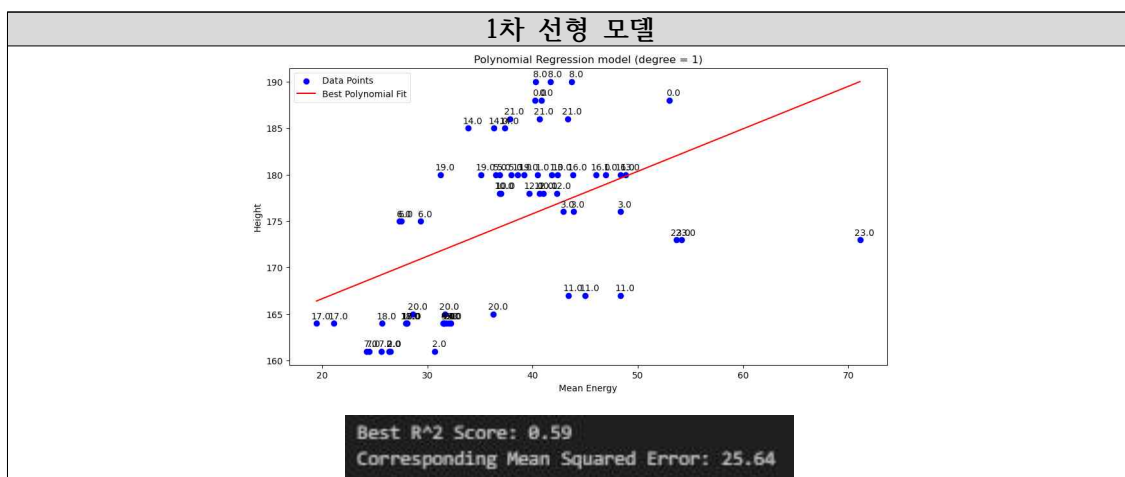
    # 기존 데이터 산점도 및 예측 곡선 그리기
    plt.figure(figsize=(12, 6))
    plt.scatter(X_scaled[:, 0], y, c='blue', label='Data Points')

    # 각 점의 id 표시
    for i, txt in enumerate(ids):
        plt.annotate(txt, (X_scaled[i, 0], y[i]), textcoords='offset points', xytext=(5, 5), ha='center')

    # 최적 모델의 예측 곡선 그리기
    X_fit = np.linspace(X_scaled.min(), X_scaled.max(), 100)[:, np.newaxis]
    y_fit = best_model.predict(X_fit)
    plt.plot(X_fit, y_fit, color='red', label='Best Polynomial Fit')

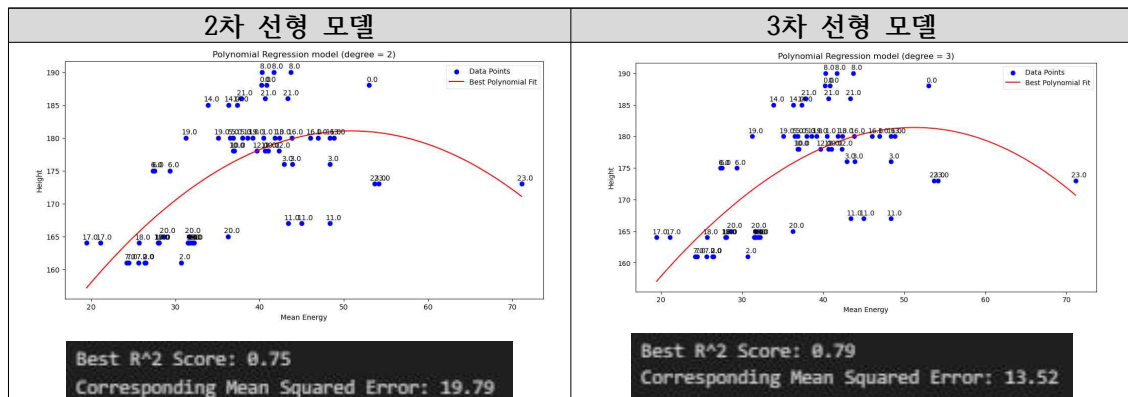
    plt.xlabel('Mean Energy')
    plt.ylabel('Height')
    plt.title(f'Polynomial Regression model (degree = {degree})')
    plt.legend()
    plt.show()
```

이후 과정에서 수월한 학습을 위해 함수를 만들어 사용하였다. test 사이즈는 20%, 그리고 반복 횟수는 1000으로 설정하였으나 입력받아 수정가능하도록 하였다. 선형 모델의 차수를 입력받고 모델을 예측한 후 결정계수와 그래프를 출력하도록 하였다. 다음은 1차 선형 모델을 사용한 결과이다.



1000의 반복을 통해 얻은 가장 높은 결정 계수를 가지는 값은 0.59로 매번 랜덤한 값을 넣기 때문에 매번 다른 R² 점수를 확인할 수 있으나 대략 0.6정도의 결정계수를 가졌다. 위 그래프에 따르면, 현재 그래프에서 결과에 큰 영향을 줄 수 있는 id가 23인 데이터에 의해 선형 모델의 기울기가 실제보다 낮게 나왔을 것으로 예상해 볼 수 있다. 이론적으로 신장에 비례하여 보폭이 넓어지고, 보폭이 넓어질수록 한 걸음당 에너지가 커질 것이므로 선형

모델을 사용하는게 적합하나, 2차, 3차 모델에 대해서도 확인해보았다.

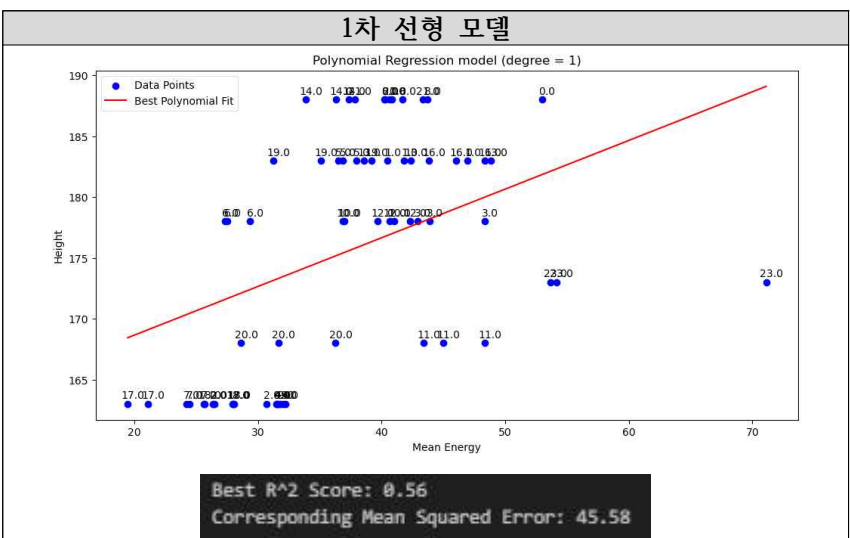


다항식을 사용한 결과 2차에서 0.75, 3차에서 0.79의 결정계수를 보였다. 이전보다 높은 결정계수를 보이지만, 이러한 모델은 경험적인 데이터와 일치하지 않으며 과적합되었을 가능성이 높음으로 1차 선형 모델을 사용하는 것이 적합할 것이다.

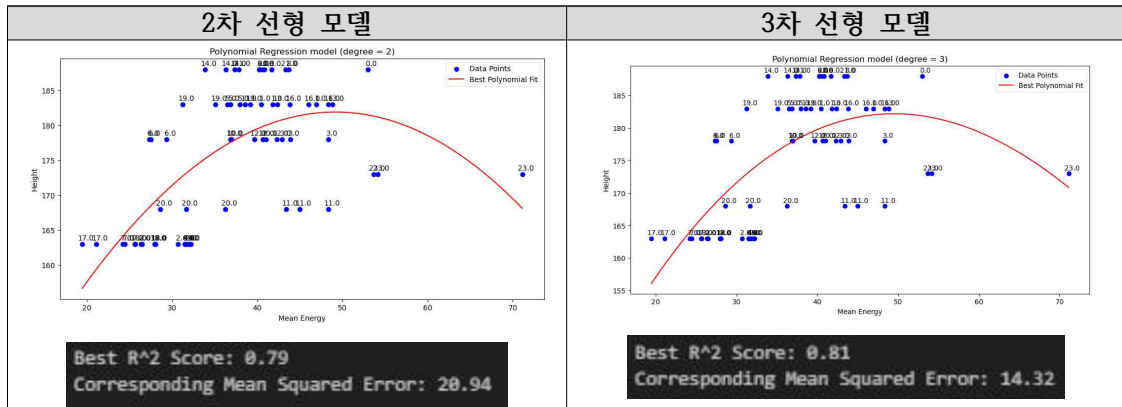
#양자화 적용

이번 과정에서는 신장에 대해 양자화를 수동으로 실시하고 모델을 학습해보았다. 우측과 같이 6단계로 나눠 선형 모델과 다항식 모델을 학습한 결과는 다음과 같다.

```
1 def quantize_height(height):
2     if height < 165:
3         return 163
4     elif 165 <= height < 170:
5         return 168
6     elif 170 <= height < 175:
7         return 173
8     elif 175 <= height < 180:
9         return 178
10    elif 180 <= height < 185:
11        return 183
12    else:
13        return 188
```



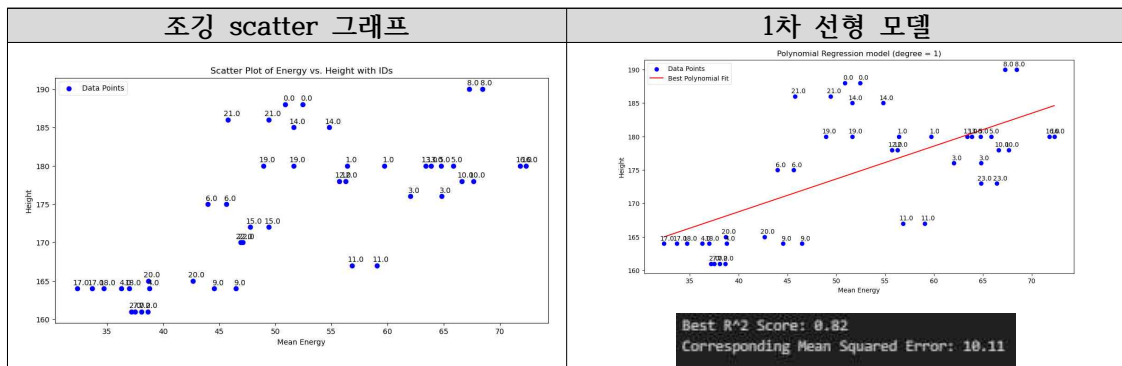
1차 선형 모델은 결정계수가 0.56로 낮아진 것을 확인할 수 있었다. 현재 사용하는 데이터셋의 샘플 수가 적음으로 양자화를 통해 정확도를 향상 시켜보고자 하였으나, 예상과는 달리 오히려 결정계수가 낮아졌다. 2차와 3차 모델에 대해서도 학습을 시켜 추가적인 확인을 해보았다.



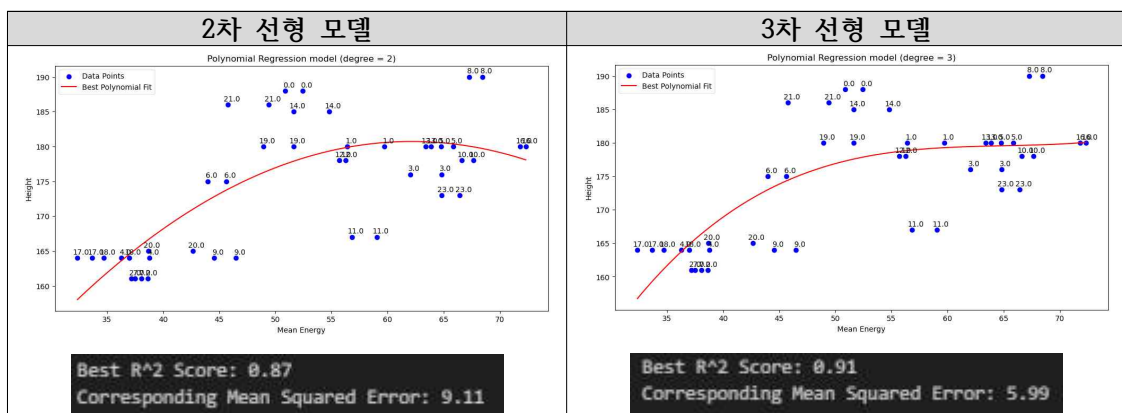
2차 선형 모델에서는 0.79, 3차에서는 0.81로 거의 양자화 이전과 큰 차이가 없는 것을 확인할 수 있었다.

결과적으로 걷기 데이터에서는 양자화를 하지 않고, 1차 선형 모델을 사용하는 것이 가장 적합할 것이다.

#모델 학습(jog)

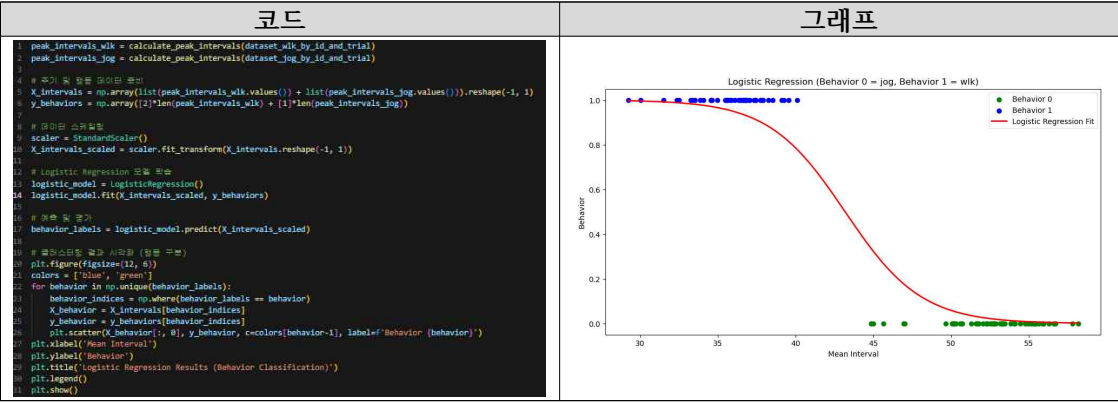


조깅 데이터의 경우 육안으로도 걷기 데이터보다 이상치가 적은 것을 확인할 수 있었다. 실제로 1차 선형 모델은 0.90으로 걷기 데이터에 비해 높은 결정계수를 보였다. 이어서, 2차, 3차 선형 모델도 확인해보았다.



2, 3차 선형 모델에서는 각각 0.87, 0.91로 1차 선형 모델보다 높은 결정계수를 보이긴 하나, 마찬가지로 과적합 되었을 가능성이 높아보인다. 따라서 조깅 데이터셋에서도 1차 선형 모델을 사용할 것이다.

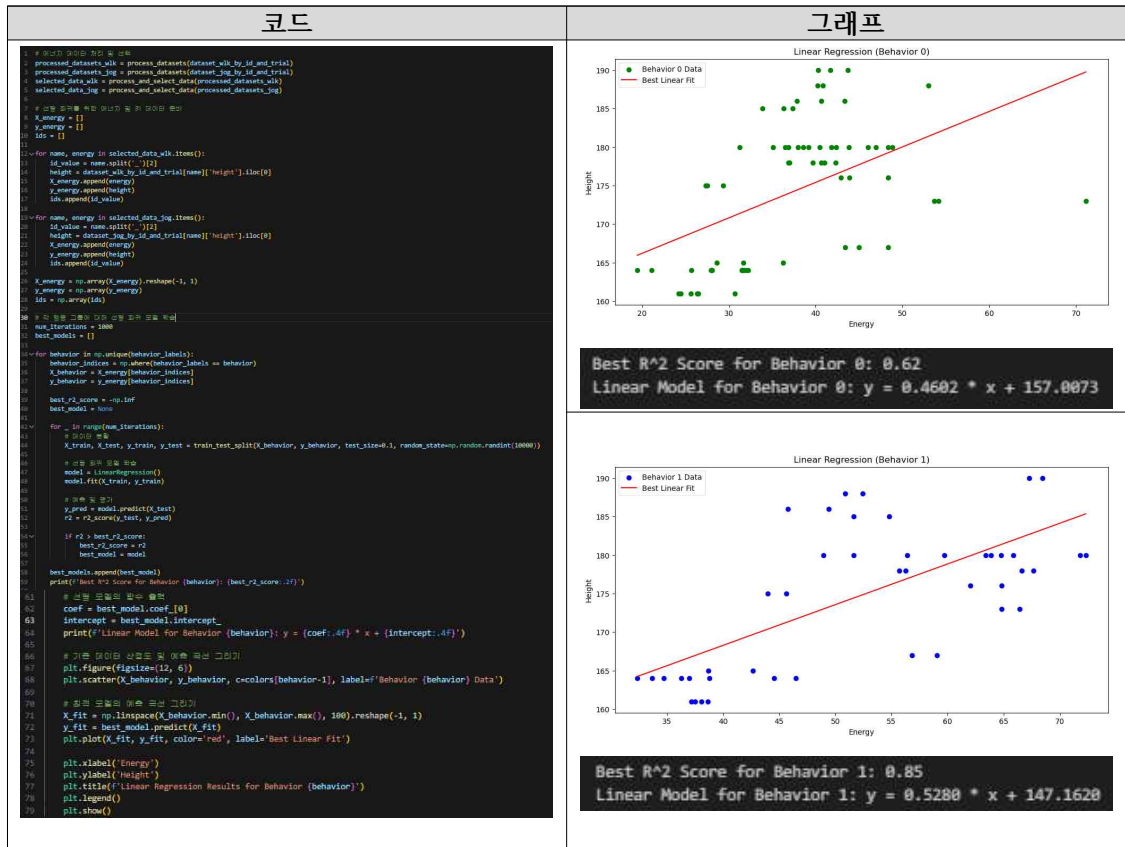
#최종 모델



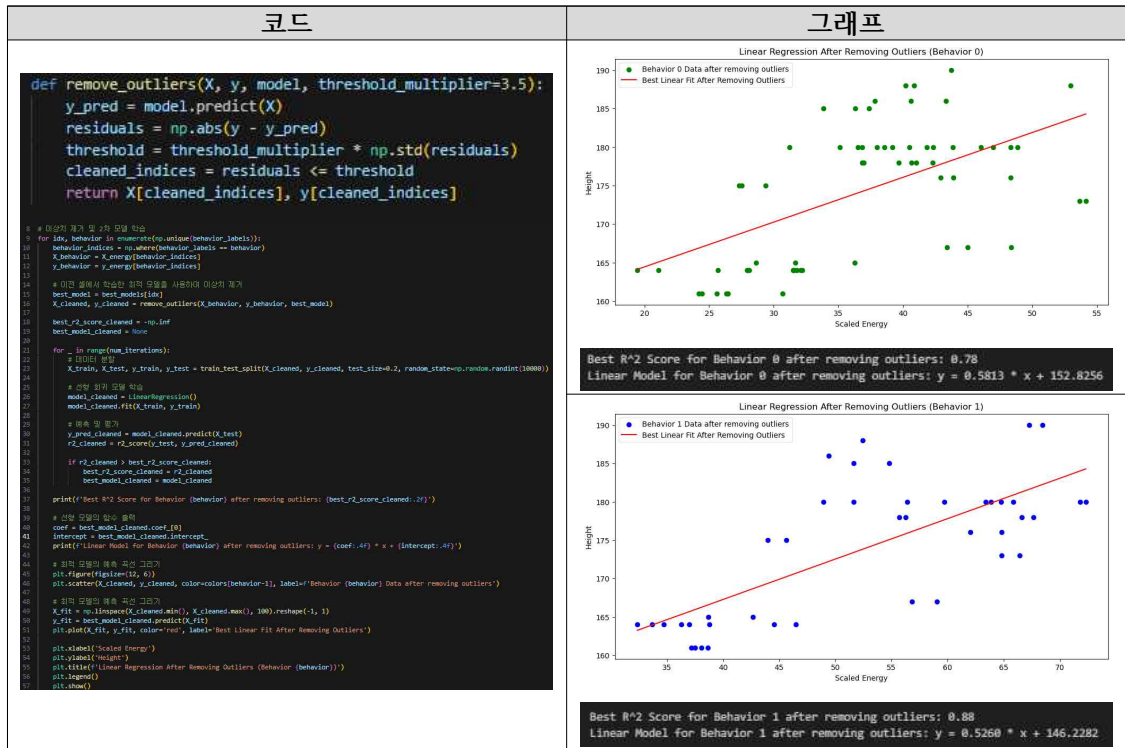
본격적인 모델 학습을 위해 주기성을 가진 데이터를 추출하여 이를 분석한다면 걷기 동작인지, 뛰기 동작인지 구분할 필요가 있다. 따라서, Logistic Regression에서 주기 당 걸음 간격을 사용하여 걷기 동작인지, 뛰기 동작인지를 먼저 구분하고, 이를 토대로 두 데이터를 구분하여 학습할 것이다.

이를 위해 wlk 데이터를 행동 0, jog 데이터를 행동 1로 구분하여 Logistic Regression을 사용하여 분류해본 결과이다. 앞서, 두 데이터의 주기 간격은 서로 겹치지 않고 분리되어있으므로 100%의 정확도로 구분된 것을 확인할 수 있다.

이후, 각 라벨로 분리하여 걷기 데이터와 조깅 데이터를 학습시켜 보았다.

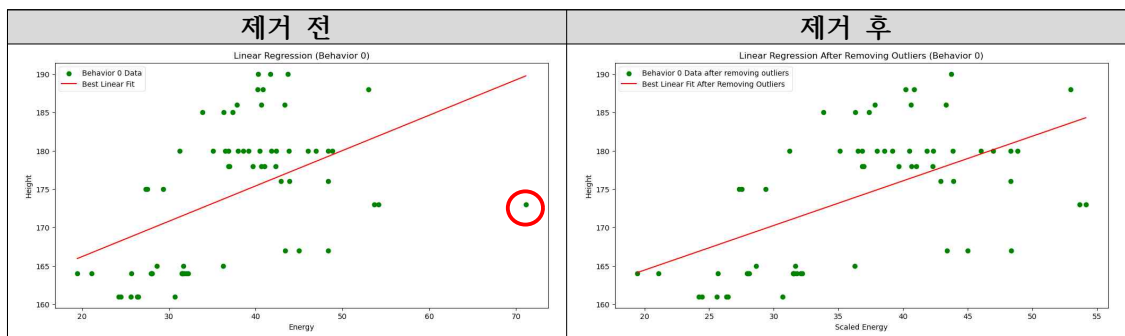


위 데이터는 행동별로 구분된 데이터를 행동에 따라 선형 학습한 결과이다. 걷기 데이터의 경우 0.74, 조깅 데이터의 경우 0.89로 이전과 이전에 개별로 확인해본 결과가 같은 것을 확인할 수 있었다. 조깅 데이터의 경우 나쁘지 않은 결과를 보였으나, 걷기 데이터의 경우 이전에 확인했던 이상치 데이터를 제거하면 더욱 좋은 모델을 얻을 수 있을 것이라 예상하였다. 따라서, 이상치를 제거하고 학습하는 2차 학습을 진행하고자 하였다. 현재는 데이터셋이 적으므로 수동으로 데이터를 제거하는 것이 가능하지만, 데이터셋이 많아지면, 수동으로 데이터를 하나씩 제거할 수 없음으로, 1차적으로 학습한 모델을 토대로 이상치를 제거하고, 2차 학습을 진행하는 방식으로 구현하였다.



2차 학습에서는 1차로 학습했던 모델을 통해 얻은 예측값과 실제값의 차이를 얻어 차이들의 표준편차를 계산하고 표준편차의 3.5배(여러번의 시도를 통해 선정하였음)를 초과하는 데이터 포인트를 이상치로 간주하고 제거한 후 재학습을 한 것이다. 결과 조깅과 걷기는 각각 결정계수 0.82와 0.91으로 걷기 데이터에서 결정계수가 향상되었으며 두 데이터셋에서 데이터셋에 대한 설명이 가능한 적절한 모델을 얻었다고 판단하였다.

실제로 이상치를 제거한 것과 그렇지 않은 두 데이터를 직접 비교하면 다음과 같다.



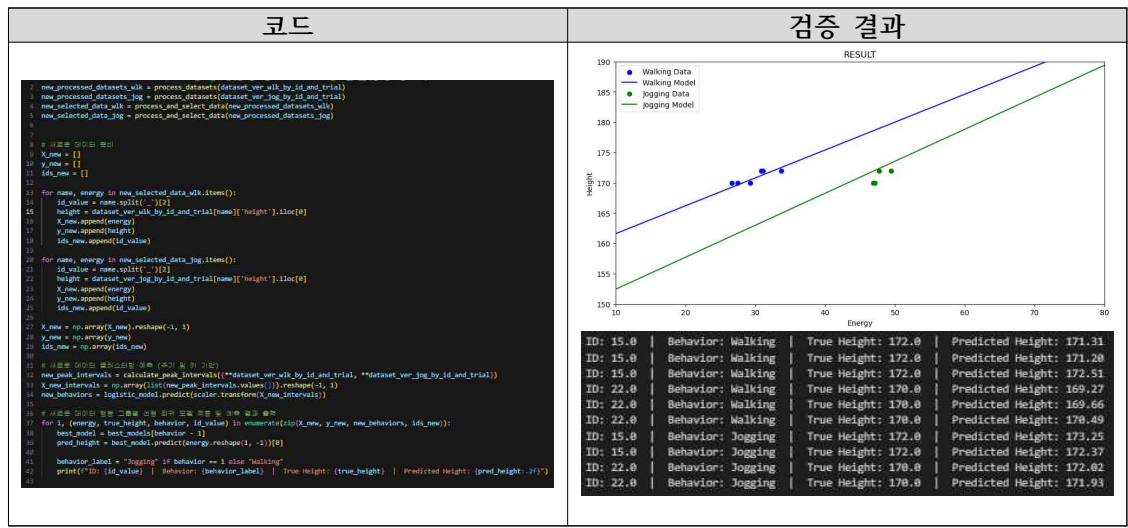
제거 전 빨간 원으로 표시한 모델 추론에 큰 방해가 될 것으로 예상하던 데이터는 필터링을 통해 일부 제거되었으며, 그 외의 데이터는 그대로 유지된 것을 확인할 수 있다. 이를 통해 2차 학습에서는 더욱 정확한 결과를 얻을 수 있을 것이다.

1차 시도와, 2차 시도에서 모델의 변화는 다음과 같다.

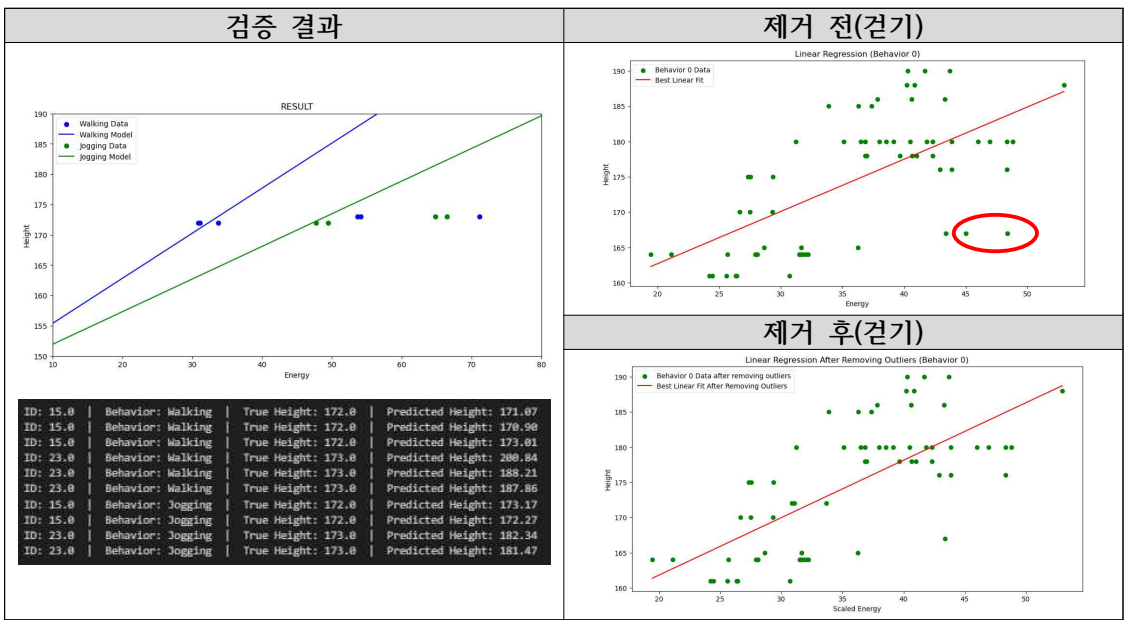
걷기	조깅
1차 : $y = 0.4602x + 157.0073$	1차 : $y = 0.5282x + 14.1620$
2차: $y = 0.5813x + 152.5256$	2차 : $y = 0.5260x + 146.2282$

조깅 데이터의 경우 1차와 2차에서의 차이가 거의 근소하나, 걷기 데이터의 경우 유의미하게 기울기가 바뀐 것을 확인할 수 있다. 필터링을 통해 정확한 모델을 얻어 사용할 것이다.

#모델 검증



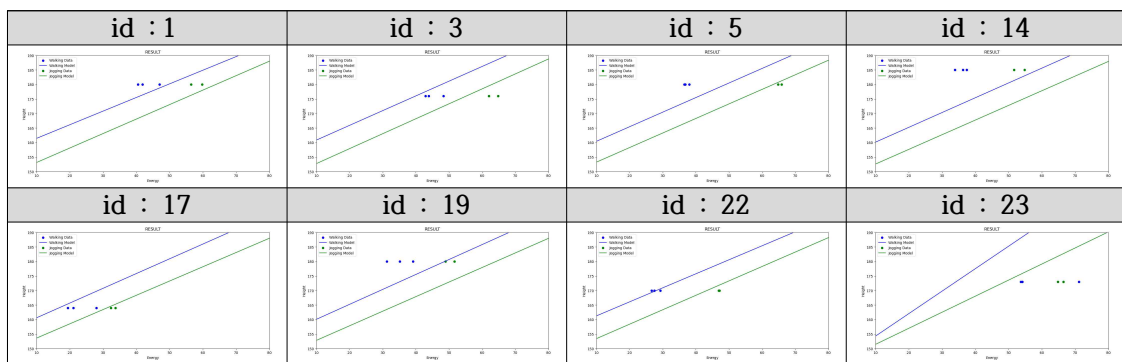
초기 분리해뒀던 데이터를 통해 추론을 진행해 보았다. 결과, id가 22와 23인 데이터 모두 높은 수준으로 키를 예측한 것을 확인할 수 있었다. 하지만, 두 데이터는 우연히 모델에 잘 맞은 것일 수 있으므로 높은 이상치를 가지던 id 23을 22 대신 검증 데이터로 하여 처음부터 다시 진행해 보았다.



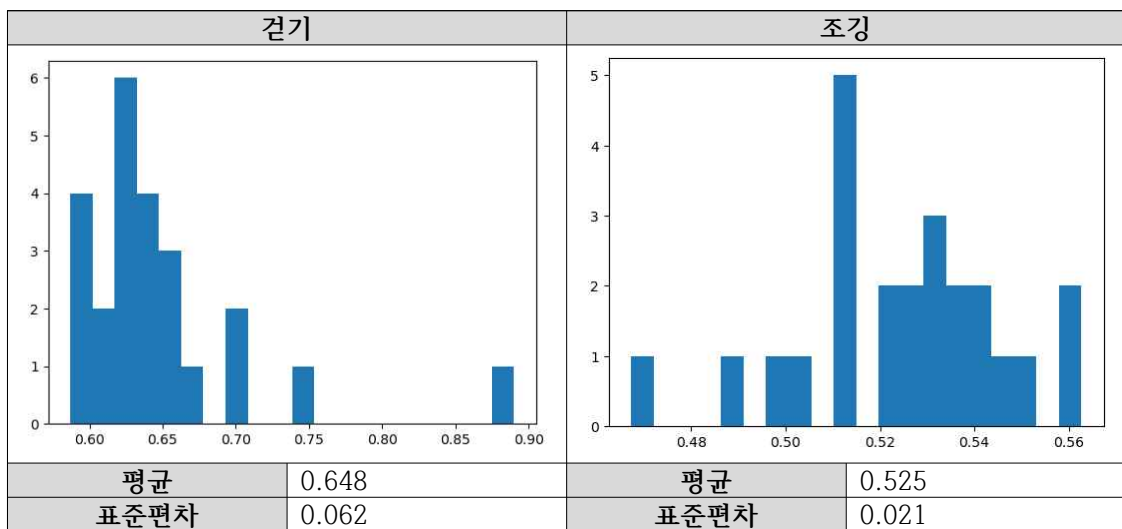
걷기	조깅
1차 : $y = 0.7388x + 147.9311$	1차 : $y = 0.5452x + 146.3989$
2차: $y = 0.8155x + 145.5199$	2차 : $y = 0.5527x + 145.2811$

이상치가 높은 데이터를 사용하니 예측 결과가 키가 173인 사람의 키를 200으로 예측하는 등, 잘못된 예측을 보여주는 것을 확인할 수 있었다. 또한, 조강의 경우 모델이 크게 달라지지 않았지만, 걷기의 경우 기울기가 크게 변한 것을 볼 수 있다. 이전, 필터링 과정에서 높은 이상치를 보이던 데이터를 필터링을 통해 제거하기 때문에, 어떤 데이터를 검증 데이터로 사용하더라도 모델에 큰 변화가 없을 것이라 예상하였다. 실제로 필터링 결과 이상치가 높은 데이터를 제거하는 것은 유효하였다. 하지만, 이상치에 대해 추론을 정확히 하지 못하는 것을 떠나, 모델 자체가 크게 달라지는 것은 추가적인 검증이 필요해보였다.

따라서, 모델 추론에 사용되는 전반적인 과정을 모아 함수로 선언하여 id가 0번인 데이터부터 23번인 데이터까지 모든 데이터에 대해 하나씩 검증 데이터로 분리하고, 모델을 학습해보았다. 밑은 그 중 일부이다.

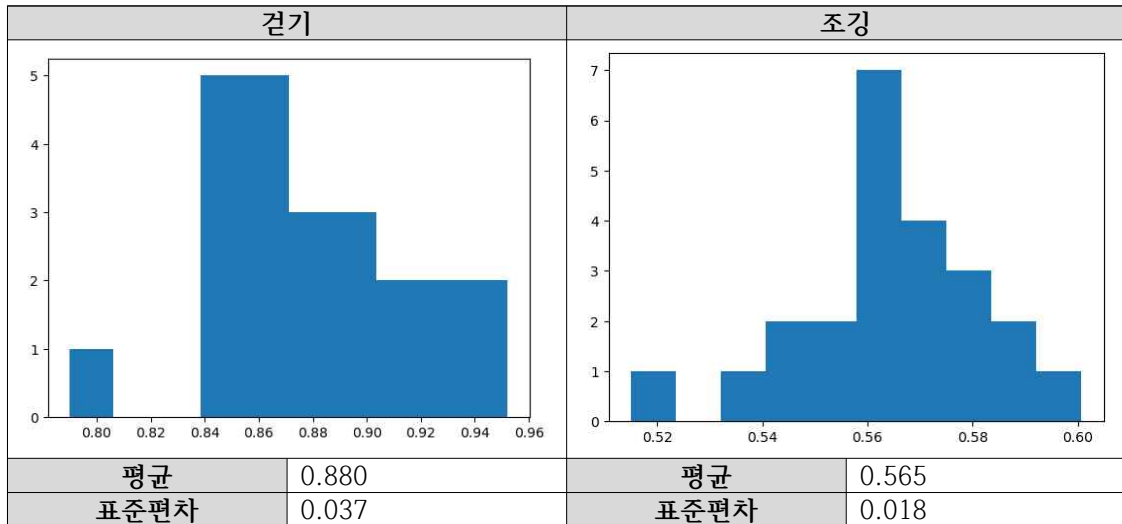


전반적으로 검증 데이터로 분리한 데이터가 얼마나 잘못된 값을 보이든 간에 비슷한 모델을 보였으나, id가 23인 모델만 다른 모델과 크게 다른 것을 볼 수 있었다. 이는 23번 데이터셋이 모델 학습에 큰 영향을 미치는 것이 분명했다. 현재, 24개의 모델의 기울기를 히스토그램으로 표출하면 다음과 같다.



걷기 데이터의 경우 특히 기울기가 0.6과 0.65 사이에 밀집되어있으며 조강 데이터의 경우도 대부분이 0.50에서 0.56 사이에 있는 것을 확인할 수 있다. 전반적으로 대부분의 기울기는 일치하나, 걷기 데이터의 경우 0.88에 해당하는 매우 높은 기울기를 가지는 데이터가 한 개 존재하는 것을 볼 수 있다. 이는 id가 23인 데이터로 마찬가지로 id 23인 데이터를 제외하면

큰 높은 기울기가 나오며 큰 영향을 준다는 것을 알 수 있다. 이를 통해 id 23인 데이터는 모델 학습에 있어서 너무 큰 영향을 줌이 확실함으로 23을 제외하고 모델을 학습하고 확인해보았다.



결과, 전체적으로 걷기 데이터는 평균이 증가하였으며, 표준편차 역시 줄어들었다. 23번 데이터가 큰 영향을 미치고 있던 것은 사실이었으며, 실제로 평균은 0.880으로 높게 나오는 것이 옳다는 것을 예상할 수 있다. 따라서 최종적으로 얻은 모델은 기울기가 평균값을 가지는 모델을 사용하는 것으로 선택하는 것이 옳을 것이다. 다음은 최종적으로 선택한 걷기와 조깅의 모델이다.

걷기	조깅
$y = 0.8809x + 143.2086$	$y = 0.5655x + 144.0638$

#결론

최종적으로 걷기와 조깅에 대한 선형 모델을 얻을 수 있었다. 하지만, 위 모델은 많은 개선 여지가 있다. 우선, 더 많은 데이터 셋을 사용하면 더욱 뛰어난 모델을 얻을 수 있을 것이다. 이전, 큰 이상치를 가지며 모델 학습에 큰 영향을 주었던 id가 23번인 데이터는 많은 데이터셋을 사용하여 평균에 가까운 데이터를 많이 수집함으로써 별도의 제거 없이 극복 가능할 것이라 예상할 수 있다. 또한, 이상치 제거, 즉 필터링 과정에서 1차적으로 학습한 모델에서 이상치를 제거하는 것이 아닌, 양자화를 통해 키를 여러 단계로 나누고, 같은 키에서 평균을 구해 이상치를 제거하면 더욱 뛰어난 모델을 얻을 수 있을 것이다. 만약, 충분한 데이터셋이 있다면, 성별에 따른 구분도 가능할 것이다. 두 번째로 센서의 민감도를 높여 더욱 정확한 센서값을 얻을 수 있다면 더욱 신뢰성 있는 에너지 값을 얻을 수 있을 것이다.

다만, 위 과정을 통해서도 해결되지 않는 문제는 있을 것이다. 보폭이라는 것이 일반적으로 신장에 따라 비례하는 것은 사실이나, 구체적인 체형, 그리고 개개인의 차이에 따라 보폭은 크게 좌지우지될 수 있어 잘못된 예측 결과를 얻을 수 있을 것이다. 이러한 한계점은 현재의 모델에서는 극복 불가능하지만, gps와 같은 다른 센서 데이터를 사용하는 등 여러 센서 데이터와 모델들을 종합하여 실제 핸드폰이 분실되었을 때 유용하게 사용할 수 있을 것이다.