



ELEC 390

Project Report

Group 33

Lauren Steel, Jacob Schaffer, Ethan Silver

20218337, 20222125, 10216012

19les5@queensu.ca, 19jsbs@queensu.ca, 19es25@queensu.ca

Friday April 14, 2023

Introduction

This report outlines the development steps taken to create a simple AI system which was then deployed in a desktop app. The process is broken down into seven steps: data collection, storage, visualization, pre-processing, feature extraction, training the classifier, and implementing the model created in a desktop app. Each phase plays a significant and necessary role in the development of a reliable and consistent AI product. Data collection is the process of collecting and measuring accurate data from several sources. Data storage is when sets are organized in a manner which is easy and convenient for use. Data visualization allows for trouble shooting with data early in the design process. Data pre-processing is the concept of fixing minor or major issues with data, to ensure that sets are usable. Feature extraction is the process of identifying and isolating various ‘features’ that contain key information, by utilizing these values the whole signal is not required. Using the features from the preprocessed training set, a logistic regression model is trained to classify the data into ‘walking’ or ‘jumping’ classes. Finally, the model could be deployed to a desktop application through use of the Tkinter library. Together these steps lead to the creation of a final AI product which can differentiate between motions of walking and jumping.

Data Collection

The first step in the development process is to collect any necessary data. Data collection is the process in which one collects and measures accurate data from multiple relevant sources. In the case of this design challenge, the goal is to develop an application that can distinguish between walking and jumping. The particular purpose the data serves is to act as representation of what a typical walking motion looks like in comparison to that of jumping. The data to be collected can be considered what is called free-living data. This type of data is recorded under real conditions without intervention to alter the environment in which the data is being recorded. In general, it is more challenging to process and extract information from a free-living dataset [1]. The reason being that they usually contain noise and are collected with a high degree of freedom . Handling these additional difficulties is addressed later in the design process with preprocessing filters.

To accurately obtain datasets representing both walking and jumping, an app called Phyphox was used to collect accelerometer data from each member’s smartphone. An accelerometer monitors acceleration in m/s^2 (or g) along three axes, x, y, and z. Figure 1 demonstrates the Phyphox interface when a 5 second accelerometer sample is produced, graphs correlating to all three axes are developed. It should be noted that an accelerometer is very sensitive to vibrations and movement, noise and artifacts are likely to be picked up [1].

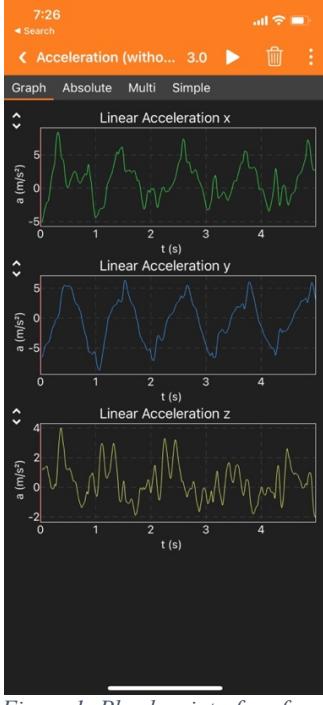


Figure 1: Phyphox interface for a 5 second accelerometer sample displaying graphs correlating to all three axes.

During the collection phase a crucial aspect to consider is diversity. A diverse dataset will allow a system to be more successful when fully developed. To maximize diversity all members of the design team collected their own set of data which would demonstrate their unique walking stride and typical jump. Moreover, each member was to ensure that throughout collection of their dataset their phone would be placed in many different positions. The variety of positions included all pant pockets, front and back, several jacket pocket samples, as well as directly carrying the phone in hand while executing the desired motions. With that being said it was ensured that the datasets collected were relatively balanced. The amount of time collected from each user, for each motion of walking or jumping, and each position would be roughly the same. The goal was to obtain a duration of data exceeding five minutes between both walking and jumping motions. The approximate outline for collection was that members would obtain thirty seconds of data for positioning in, front pockets, back pockets, jacket pockets, and being held in hand for both jumping and walking samples. Resulting in approximately two and a half minutes for each walking and jumping.

Data samples collected from the app Phyphox could be directly exported as CSV files. Therefore, data could be easily moved in any fashion most convenient to the user (email, airdrop, etc.). The transferred zip files included the CSV along with correlating meta data. Once these files for all collected samples were transferred to a member's PC it was possible to label the data as fit. Each file specified the movement (walking or jumping), position (front pocket, back pocket, jacket pocket or hand) and side (left or right) that the corresponding data represented.

Two members collected the data in thirty second segments. One member chose to collect data directly in the necessary five second samples. Those who collected samples of longer segments took different approaches in splitting the data. One member manually broke up the CSVs into

five second periods according to time stamps. The other approach used was to directly add the larger CSV files containing thirty seconds of walking or jumping acceleration data into a python file. Within the program it is possible to break up the data as shown in Figure 2.

```
Q5 = jump5
B5 = jump5.iloc[:, 0]
H51 = Q5[B5 < 5]
H52 = Q5[(5 < B5) & (B5 < 10)]
H53 = Q5[(10 < B5) & (B5 < 15)]
H54 = Q5[(15 < B5) & (B5 < 20)]
H55 = Q5[(20 < B5) & (B5 < 25)]
H56 = Q5[25 <= B5]

with h5py.File('./data.h5', 'w') as hdf:
    G1 = hdf.create_group('/Ethan')
    W = G1.create_group('/Walk')
    W.create_dataset('walk1-1', data=Z11)
    W.create_dataset('walk1-2', data=Z12)
    W.create_dataset('walk1-3', data=Z13)
    W.create_dataset('walk1-4', data=Z14)
    W.create_dataset('walk1-5', data=Z15)
    W.create_dataset('walk1-6', data=Z16)
    W.create_dataset('walk2-1', data=Z21)
```

Figure 2: A segment illustrative of the CSV to HDF5 conversion program.

Data Storage

Once data has been collected, it must be stored. Data storage should be organized in a manner which is easy and convenient for usage. After transferring all the subsets of data to a personal computer and adding labels, each member stored their collected data in the form of a HDF5 file. HDF5 is a data model as well as a library, and file format for the storage and management of data. HDF5 files are designed to handle large, complex data with an open format, meaning they are supported by many tools such as Python [2].

An HDF5 file can be viewed as a container or grouping, which holds datasets. The datasets could be a range of various file types. The two important objects in HDF5 files are ‘groups’ and ‘datasets’ [2]. Groups are like folders implemented for organization. Every HDF5 file includes a group called ‘root’. Dealing with groups and their content is very similar to path directories. A ‘/’ represents the root group and different members will follow the backslash to branch off into subgroups.

The complete HDF5 file was organized in a tree with four main groups. Three of these groups would correspond to each member’s personal datasets. The final group of the requested file is labeled ‘dataset’ and has two subgroup members corresponding to either train or test data.

Creating train and test splits is vital to create a simple AI system [2]. After the team ensured each signal of their personal data was in the form of five second windows, as described in the data collection section, the data was then shuffled. 90% of collected data is used for training purposes and the remaining 10% is added to the testing subgroup. So, within the ‘dataset’ group, the contents of each subgroup train and test were taken directly from all the individual’s personally collected files but reorganized and stored with a different structure. It is typically good practice to keep the data as originally collected, which is why the groups corresponding to team members remain in the final HDF5 file. A visual of the final data organization within the HDF5 file is represented in Figure 3.

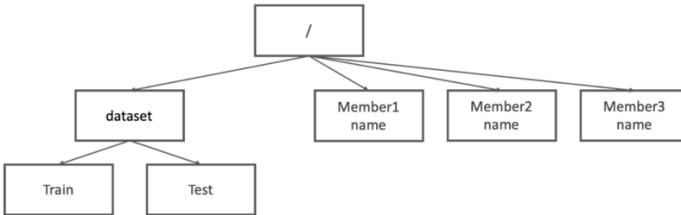


Figure 3: Final HDF5 file hierarchy visualization, Image from [3].

Each individual team member created their own HDF5 file to store their data in which they were able to organize sets as they desired. These files would likely take the form of two groups ‘walking’ and ‘jumping’ which could then branch down to various subgroups of phone positions collected (pant, jacket, hand) and the side it was taken (right, left). To start, we first include the necessary library h5py, then begin creating and writing to the HDF5 file. To ensure the proper organization it is necessary to use ‘hdf.create_group’ within python. This allows us to begin a hierarchy of groups.

To replicate the desired structure seen above, a separate python file is created to read individual team members HDF5 files into one large hierarchy. The implementation code written to accomplish this can be seen in Figure 4 and Figure 5. Where the files ‘ethan.h5’, ‘Jacob_Schaffer.h5’ and ‘lauren_data.h5’ correspond to individual’s collected data which are represented as Member1 name, Member2 name etc. in Figure 3.

```

1 import h5py
2
3 ethan = h5py.File('ethan.h5', 'r')
4 jacob = h5py.File('Jacob_Schaffer.h5', 'r')
5 lauren = h5py.File('lauren_data.h5', 'r')
6 traintest = h5py.File('train_test.h5', 'r')
7 dataout = h5py.File('data.h5', 'w')
8
9 dataout.create_group('Ethan')
10 GEJ = dataout.create_group('Ethan/Jump')
11 GEW = dataout.create_group('Ethan/Walk')
12 dataout.create_group('Lauren')
13 GLJ = dataout.create_group('Lauren/Jump')
14 GLW = dataout.create_group('Lauren/Walk')
15 dataout.create_group('Jacob')
16 GJJ = dataout.create_group('Jacob/Jump')
17 GJW = dataout.create_group('Jacob/Walk')
18 dataout.create_group('dataset')
19 GTrain = dataout.create_group('dataset/Train')
20 GTest = dataout.create_group('dataset/Test')
21
22 # Add Ethan Data
23 for dataset in ethan['Jump'].keys():
24     GEJ.create_dataset(name=dataset, data=ethan['Jump'][dataset])
25 for dataset in ethan['Walk'].keys():
26     GEW.create_dataset(name=dataset, data=ethan['Walk'][dataset])
27
28 # Add Lauren Data
29 for dataset in lauren.keys():
30     if 'jumping' in dataset:
31         GLJ.create_dataset(name=dataset, data=lauren[dataset])
32     if 'walking' in dataset:
33         GLW.create_dataset(name=dataset, data=lauren[dataset])
34

```

Figure 4: Part 1 of Python code from h5_generator.py which creates an HDF5 file with desired organization.

```

35  # Add Jacob Data
36  for dataset in jacob['Jumping'].keys():
37      GJJ.create_dataset(name=dataset, data=jacob['Jumping'][dataset])
38  for dataset in jacob['Walking'].keys():
39      GJW.create_dataset(name=dataset, data=jacob['Walking'][dataset])
40
41 # Add Train Data
42 for dataset in traintest['dataset']['Train'].keys():
43     GTrain.create_dataset(name=dataset, data=traintest['dataset']['Train'][dataset])
44
45 # Add Test Data
46 for dataset in traintest['dataset']['Test'].keys():
47     GTest.create_dataset(name=dataset, data=traintest['dataset']['Test'][dataset])

```

Figure 5: Part 2 of Python code from file `h5_generator.py` which creates an HDF5 file with desired organization.

Visualization

Data visualization allows for trouble shooting with data early in the design process. It is a critical step in the field of data science, and it is a means of finding issues as well as becoming more familiar with the data. Overall, visualization helps to gain insights, identify patterns, trends, outliers, and summarize information in bigger data sets [4]. In this stage, it is necessary to visualize a few samples from the datasets on all three axes and from both motions of walking and jumping. Data visualization will be achieved using Python library Matplotlib. As shown in the plots below, there are many features that create a distinction between walking and jumping data, and these differences are taken advantage of during feature extraction.

Handheld Data

The below figures illustrate the difference in results between jumping and walking with the phone held in the user's hand. As the figures below clearly demonstrate, the data for jumping and walking are very distinct. Jumping data has much higher positive extremes on all axes, and much lower negative extremes on the Y axis. The data also simply deviates much more, which can be utilized during feature extraction.

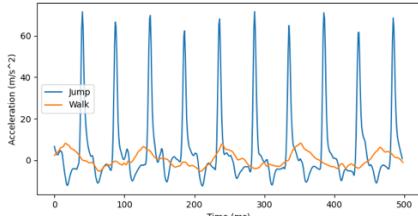


Figure 6: Handheld Acceleration data for the X axis.

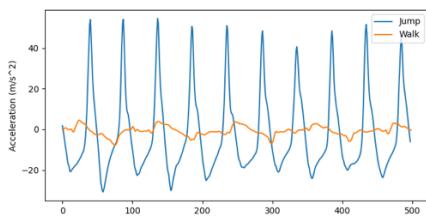


Figure 7: Handheld Acceleration data for the Y axis.

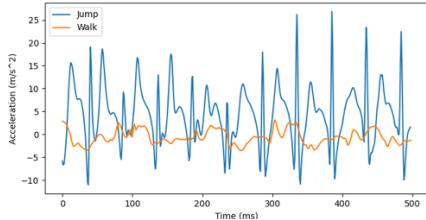


Figure 8: Handheld Acceleration data for the Z axis.

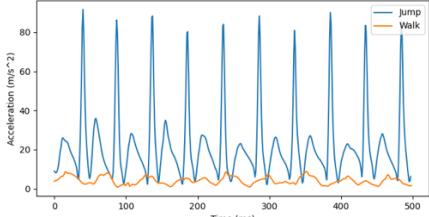


Figure 9: Handheld Absolute Acceleration data.

Pocket Data

Below are figures which contrast the accelerometer data collected while the user has their phone in their pocket. The data is much noisier and less consistent than with the phone in the user's hand, which makes sense since the phone will naturally shake around more in a pocket than in a hand. While the distinct features of the data are less obvious than with the phone in the user's hand, jumping data still has distinct features from walking data, particularly on the Y axis, as shown in Figure 11. Acceleration values reach much more extreme negatives for walking than for jumping. Another field that still expresses very distinct values is the absolute acceleration data, shown in Figure 13.

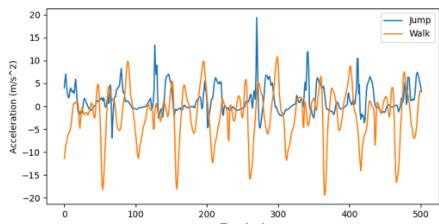


Figure 10: Pocket Acceleration data for the X axis.

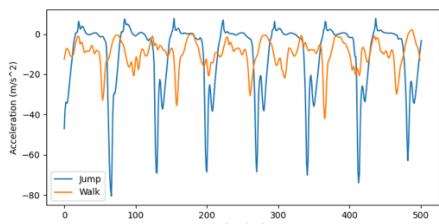


Figure 11: Pocket Acceleration data for the Y axis.

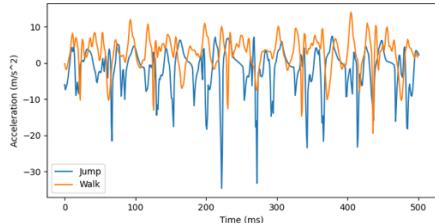


Figure 12: Pocket Acceleration data for the Z axis.

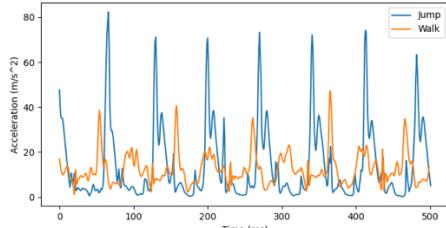


Figure 13: Absolute Pocket Acceleration data.

Future Considerations

While most of the data collected turned out clean as shown above, as it was collected in a controlled environment, real-world data may have significantly more noise. The following data was some of the noisier data collected, and if the model were to be redeveloped with more robust functionality, the variety of real-world results compared to controlled data collection would be a more important consideration.

As shown below, the differences between jumping and walking acceleration are much less pronounced, though there are still distinct features that allow it to be classified. The unique features selected and extracted will be discussed further in Pre-processing.

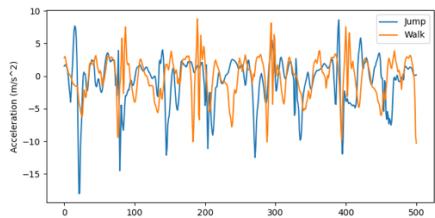


Figure 14: Real-world Pocket Acceleration for the X axis.

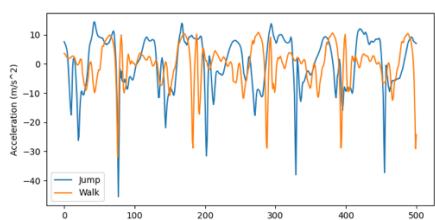


Figure 15: Real-world Pocket Acceleration for the Y axis.

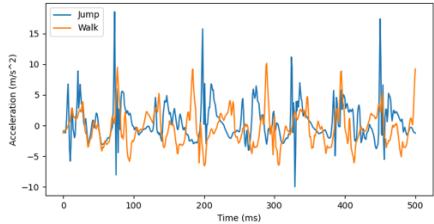


Figure 16: Real-world Pocket Acceleration for the Z Axis.

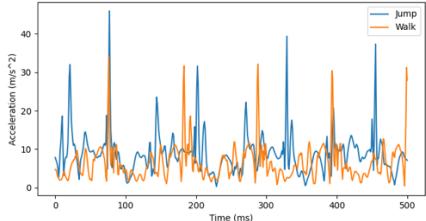


Figure 17: Absolute Real-world Pocket Acceleration.

Pre-processing

In almost every dataset, even when making extreme efforts to be careful during the data collection phase, there will still inevitably be noise. Data pre-processing is the concept of fixing minor or major issues with data, to ensure that sets are usable. This process could involve a variation of numerous approaches inclusive of imputation, normalization, feature extraction, dimensionality reduction, noise reduction, source separation, etc. [5].

In the case of this model a few of these pre-processing techniques were taken. The data collected contained plenty of noise. It is vital to minimize noise as it can destroy features in a signal. There are two categories, high-frequency noise, and low-frequency noise. In this case the noise obscuring the data is high frequency noise which adds fast moving noise to the signal. To reduce this type of noise a moving average (MA) filter is applied. MA is a low-pass filter, meaning it passes through low frequencies and discards high frequencies [5].

The idea of a MA filter is that they slide a window of a given size down the signal and average all the data points within that window. The filter will replace each value of the signal with the generated average value of values in a filter which precedes and includes said value [5]. The logic behind averaging the numbers is that the noise values cancel each other out.

The next approach was taken after the feature extraction. At this point an effort was made to detect and remove any outliers from the collected data. It was determined that, after analysing the de-noised data, none of the 5-second segments were far enough separate from each other to warrant treating those segments as outliers. As a result, the feature extraction and classification moved forward with the entire dataset. During the classification stage, trials were run with several subsets of the dataset to determine if any data was creating outliers, but the model's resulting accuracy was unaffected, reinforcing the lack of outliers in the normalized dataset.

Figure 18 below shows a graph of total acceleration vs time for a sample of walking data (blue) and a sample of jumping data (orange), Figure 19 shows the same sets of data graphed after a

moving average filter has been applied to the data. Viewing these graphs, it is clear that applying the moving average filter greatly reduced the noise in the data, removing many of the false peaks in acceleration, which results in a significantly more accurate model.

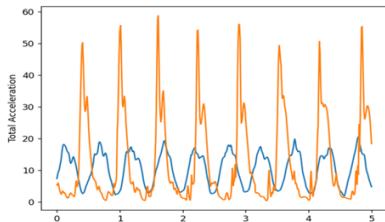


Figure 18: Absolute acceleration for a sample walking (blue) vs jumping (orange) data subset.

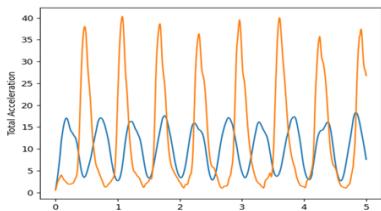


Figure 19: The same data as above, with a moving average applied to reduce noise.

Lastly, it was necessary to normalize the data so that it is fit for logistic regression. Code demonstrating these steps is displayed in the classifier program.

Feature Extraction

Feature extraction is the process of identifying and isolating various ‘features’ that contain key information. By utilizing these values, the whole signal is not required. During feature extraction there are a few design choices to be made, the first being window size. With a smaller window more information is extracted. The time window examined in this case is five second segments. Next there is window overlap. The model presented uses a rolling window. A rolling window moves ahead one sample at a time. This assists in retaining important information. Lastly and most importantly we must choose the types of features for extraction.

For this model, ten features were selected for extraction. The chosen features include, average x acceleration, maximum x acceleration, average y acceleration, maximum y acceleration, average z acceleration, maximum z acceleration, average time above the average z acceleration, overall average acceleration, overall maximum acceleration, and average time above the overall average acceleration. The goal in picking features is choosing a variety to summarize the important information contained in the original data sample [6].

Feature extraction addresses the issue of attaining the most informative and compact set of data we can, to improve the performance of a machine learning model. First, we consider ‘informative’. We are looking for features that can characterize the behavior of what we are trying to model [7]. In this case the end goal of the model is to classify a given data sample as either ‘walking’ or ‘jumping’. So, the strategy taken is to extract features that would be the most useful in defining a signal into one of those two categories.

What is meant by 'compact', is that we want to exclude irrelevant features from the model [7]. This proved to be an important consideration in the simple AI implementation. Throughout the design process several combinations of features were tested. It was discovered through trial and error that with too many feature extractions the model was less accurate in comparison to tests with fewer features. Further analysis would be necessary to select the distinct features of which were most beneficial in generating frequently successful results. Upon analyzing data visualizations of both jumping and walking samples and comparing the two, conclusive observations were made. It was determined that jumping data yielded acceleration vs. time graphs with overall higher and shorter peaks than those of walking samples. Therefore, a focus was set upon the average and maximum accelerations of each sample across all axes as well as averaged overall. This focus caused accuracy to improve drastically. A final ten features which yielded the highest accuracy were determined. Code implementing the feature extraction is shown below as Figure 20 to Figure 22. Note that the script shown below is extraction from the HDF5 file data.h5 from the group Ethan/Jump, the rest of this python file includes extraction from groups pertaining to every members data of both jumping and walking samples. Figure 23 represents the last few lines of code in the file writing to project_files.h5 in groups train and test.

```

1 import h5py as h5
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import pandas as pd
5 from sklearn.model_selection import train_test_split
6
7 empty = []
8 emptyset = np.array(empty)
9
10 #csv_file = "Feature_Extraction_Data.csv"
11 #info = ["avgXAccel", "maxXAccel", "avgYAccel", "maxYAccel", "avgZAccel", "maxZAccel", "avgTimeAboveAvgZ", "avgAccel", "maxAccel", "avgTimeAboveAvg", "label"]
12
13 #df.to_csv(csv_file, mode='a', header=False, index=False)
14
15
16 df = pd.DataFrame([emptyset])
17
18 with h5.File('../data.h5', 'r') as hdf:
19     group = hdf.get('Ethan/Jump')
20     group = list(group.keys())
21     for i in range(1, len(group)+1):
22         with h5.File('../data.h5', 'r') as hdf:
23             dataset = hdf.get('Ethan/Jump/' + str(group[i - 1]))
24             myarray = np.array(dataset)
25             #print(dataset)
26             for i in range(0, len(myarray)):
27                 if i == 0:
28                     myarray[i][1] = myarray[i][1]
29                     myarray[i][2] = myarray[i][2]
30                     myarray[i][3] = myarray[i][3]
31                 else:
32                     myarray[i][1] = (myarray[i][1] + myarray[i - 1][1] + myarray[i - 2][1] + myarray[i - 3][1]) / 4
33                     myarray[i][2] = (myarray[i][2] + myarray[i - 1][2] + myarray[i - 2][2] + myarray[i - 3][2]) / 4
34                     myarray[i][3] = (myarray[i][3] + myarray[i - 1][3] + myarray[i - 2][3] + myarray[i - 3][3]) / 4
35
36             # Average total acceleration
37             avgAccel = 0
38             for i in range(0, len(myarray)):
39                 avgAccel += abs(myarray[i][4]) / 500
40
41             # max total acceleration

```

Figure 20: Part 1 of Python code from file Feature_Extraction_And_Storage.py which creates performs feature extraction and creates an HDF5 file called project_files.h5.

```

42     maxAccel = 0
43     for i in range(0, len(myarray)):
44         if abs(myarray[i][4]) > maxAccel:
45             maxAccel = abs(myarray[i][4])
46
47     # Average X acceleration
48     avgXAccel = 0
49     for i in range(0, len(myarray)):
50         avgXAccel += myarray[i][1] / 500
51
52     # Average Y acceleration
53     avgYAccel = 0
54     for i in range(0, len(myarray)):
55         avgYAccel += myarray[i][2] / 500
56
57     # Average Z acceleration
58     avgZAccel = 0
59     for i in range(0, len(myarray)):
60         avgZAccel += abs(myarray[i][3] / 500)
61     # average time above avg Z acceleration
62
63     avgZTime = 0
64     timer = False
65     times = 0
66     totalTime = 0
67     for i in range(0, len(myarray)):
68         if abs(myarray[i][3]) >= avgZAccel:
69             if timer == False:
70                 timer = True
71                 times += 1
72                 totalTime += 1
73             else:
74                 totalTime += 1
75         else:
76             timer = False
77
78     avgZTime = totalTime / times
79
80     # max Z acceleration
81     maxZAccel = 0
82     for i in range(0, len(myarray)):

```

Figure 21: Part 2 of Python code from file *Feature_Extraction_And_Storage.py* which creates performs feature extraction and creates an HDF5 file called *project_files.h5*.

```

83         if abs(myarray[i][3]) > maxZAccel:
84             maxZAccel = abs(myarray[i][3])
85
86     # max X acceleration
87     maxXAccel = 0
88     for i in range(0, len(myarray)):
89         if abs(myarray[i][1]) > maxXAccel:
90             maxXAccel = abs(myarray[i][1])
91
92     # max Y acceleration
93     maxYAccel = 0
94     for i in range(0, len(myarray)):
95         if abs(myarray[i][2]) > maxYAccel:
96             maxYAccel = abs(myarray[i][2])
97
98     # average time above avg acceleration
99
100    avgTime = 0
101    timer = False
102    times = 0
103    totalTime = 0
104    for i in range(0, len(myarray)):
105        if abs(myarray[i][4]) >= abs(avgAccel):
106            if timer == False:
107                timer = True
108                times += 1
109                totalTime += 1
110            else:
111                totalTime += 1
112        else:
113            timer = False
114
115    avgTime = totalTime / times
116
117
118    info = [avgXAccel, maxXAccel, avgYAccel, maxYAccel, avgZAccel, maxZAccel, avgZTime, avgAccel, maxAccel, avgTime, 1]
119    df = pd.concat([df, pd.DataFrame([info])])

```

Figure 22: Part 3 of Python code from file *Feature_Extraction_And_Storage.py* which creates performs feature extraction and creates an HDF5 file called *project_files.h5*.

```

df.columns = ["avgXAccel", "maxXAccel", "avgYAccel", "maxYAccel", "avgZAccel", "maxZAccel", "avgTimeAboveAvgZ", "avgAccel", "maxAccel", "avgTimeAboveAvg", "label"]

Train, Test = train_test_split(df, test_size=0.1, shuffle=True)

with h5.File("./project_files.h5", 'w') as hdf:
    group = hdf.create_group('dataset')
    group.create_dataset('train', data = Train)
    group.create_dataset('test', data = Test)

```

Figure 23: Part 4 of Python code from file *Feature_Extraction_And_Storage.py* which creates performs feature extraction and creates an HDF5 file called *project_files.h5*.

Training the Classifier

Using the features from the preprocessed training set, a logistic regression model was trained to classify the data into ‘walking’ and ‘jumping’ classes. During the training phase it was important to ensure that the test set didn’t leak into the training set. This implies that there was no overlap

between the segments used for training and those used for testing. Below are figures displaying the implementation of the classifier in python.

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score, recall_score, confusion_matrix, ConfusionMatrixDisplay, roc_curve, \
7     RocCurveDisplay, roc_auc_score
8 from sklearn.pipeline import make_pipeline
9 from sklearn.linear_model import LogisticRegression
10 from sklearn.inspection import DecisionBoundaryDisplay
11 from sklearn.decomposition import PCA
12 import h5py as h5
13
14
15 def Classify(model, scaler, df):
16
17     # delete first row
18     df = df.drop(df.index[0])
19     df = df.dropna()
20     myarray = np.array(df)
21
22
23
24     for i in range(3, len(myarray)):
25         myarray[i][1] = (myarray[i][1] + myarray[i - 1][1] + myarray[i - 2][1] + myarray[i - 3][1]) / 4
26         myarray[i][2] = (myarray[i][2] + myarray[i - 1][2] + myarray[i - 2][2] + myarray[i - 3][2]) / 4
27         myarray[i][3] = (myarray[i][3] + myarray[i - 1][3] + myarray[i - 2][3] + myarray[i - 3][3]) / 4
28
29     # Average total acceleration
30     avgAccel = 0
31     for i in range(0, len(myarray)):
32         avgAccel += abs(myarray[i][4]) / 500
33
34     # max total acceleration
35     maxAccel = 0
36     for i in range(0, len(myarray)):
37         if abs(myarray[i][4]) > maxAccel:
38             maxAccel = abs(myarray[i][4])
39
40     # Average X acceleration
41     avgXAccel = 0

```

Figure 24: Part 1 of Python code from file *CLASSIFIER_FINAL.py* which creates performs training on the classifier.

```

42     for i in range(0, len(myarray)):
43         avgAccel += myarray[i][1] / 500
44
45     # Average Y acceleration
46     avgYAccel = 0
47     for i in range(0, len(myarray)):
48         avgAccel += myarray[i][2] / 500
49
50     # Average Z acceleration
51     avgZAccel = 0
52     for i in range(0, len(myarray)):
53         avgZAccel += abs(myarray[i][3]) / 500
54     # average time above avg Z acceleration
55
56     avgZTime = 0
57     timer = False
58     times = 0
59     totalTime = 0
60     for i in range(0, len(myarray)):
61         if abs(myarray[i][3]) >= abs(avgZAccel):
62             if timer == False:
63                 timer = True
64                 times += 1
65                 totalTime += 1
66             else:
67                 totalTime += 1
68         else:
69             timer = False
70
71     if times > 0:
72         avgZTime = totalTime / times
73     else:
74         avgZTime = 0
75
76     # max Z acceleration
77     maxZAccel = 0
78     for i in range(0, len(myarray)):
79         if abs(myarray[i][3]) > maxZAccel:
80             maxZAccel = abs(myarray[i][3])
81
82     # max X acceleration

```

Figure 25: Part 2 of Python code from file *CLASSIFIER_FINAL.py* which creates performs training on the classifier.

```

83     maxXAccel = 0
84     for i in range(0, len(myarray)):
85         if abs(myarray[i][1]) > maxXAccel:
86             maxXAccel = abs(myarray[i][1])
87
88     # max Y acceleration
89     maxXAccel = 0
90     for i in range(0, len(myarray)):
91         if abs(myarray[i][2]) > maxXAccel:
92             maxXAccel = abs(myarray[i][2])
93
94     # average time above avg acceleration
95
96     avgTime = 0
97     timer = False
98     times = 0
99     totalTime = 0
100    for i in range(0, len(myarray)):
101        if abs(myarray[i][4]) >= abs(avgAccel):
102            if timer == False:
103                timer = True
104                times += 1
105                totalTime += 1
106            else:
107                totalTime += 1
108        else:
109            timer = False
110
111    if times > 0:
112        avgTime = totalTime / times
113    else:
114        avgTime = 0
115
116    info = [avgXAccel, maxXAccel, avgYAccel, maxYAccel, avgZAccel, maxZAccel, avgZTime, avgAccel, maxXAccel, avgTime]
117    df = pd.DataFrame(info)
118    Features = df
119    df.columns = ["avgXAccel", "maxXAccel", "avgYAccel", "maxYAccel", "avgZAccel", "maxZAccel", "avgTimeAboveAvgZ",
120                  "avgAccel", "maxAccel", "avgTimeAboveAvg"]
121
122    df = scaler.transform(df)
123    pred = pd.DataFrame(model.predict(df))

```

Figure 26: Part 3 of Python code from file CLASSIFIER_FINAL.py which creates performs training on the classifier.

```

124     final = pd.concat([Features, pred], axis=1)
125     final.columns = ["avgXAccel", "maxXAccel", "avgYAccel", "maxYAccel", "avgZAccel", "maxZAccel", "avgTimeAboveAvgZ",
126                       "avgAccel", "maxAccel", "avgTimeAboveAvg", 'label']
127
128
129
130
131
132
133 def train():
134     with h5.File('./project_files.h5', 'r') as hdf:
135         Train = pd.DataFrame(hdf.get('/dataset/train'))
136         Test = pd.DataFrame(hdf.get('/dataset/test'))
137
138     Train = Train.dropna()
139     Test = Test.dropna()
140
141     Train.columns = ["avgXAccel", "maxXAccel", "avgYAccel", "maxYAccel", "avgZAccel", "maxZAccel", "avgTimeAboveAvgZ",
142                      "avgAccel", "maxAccel", "avgTimeAboveAvg", "label"]
143     Test.columns = ["avgXAccel", "maxXAccel", "avgYAccel", "maxYAccel", "avgZAccel", "maxZAccel", "avgTimeAboveAvgZ",
144                      "avgAccel", "maxAccel", "avgTimeAboveAvg", "label"]
145
146     Y_train = Train['label']
147     Y_test = Test['label']
148
149     X_train = Train.loc[:, :-1]
150     X_test = Test.loc[:, :-1]
151
152     scaler = StandardScaler()
153     X_train = scaler.fit_transform(X_train)
154     X_test = scaler.transform(X_test)
155
156     # Now, similar to the code presented in the lecture slides, train and test your model.
157
158     l_reg = LogisticRegression(max_iter=10000)
159     clf = make_pipeline(StandardScaler(), l_reg)
160
161     clf.fit(X_train, Y_train)
162
163     return clf, scaler
164

```

Figure 27: Part 4 of Python code from file CLASSIFIER_FINAL.py which creates performs training on the classifier.

The classifier shown in the figures above includes two main functions, classify and train. Classify performs some preprocessing and feature extraction on the input df, which contains accelerometer data. Any rows with missing or irrelevant values are dropped. After calculating the features based on given accelerometer readings, they are then combined into a new DataFrame. The new DataFrame ‘Features’ is transformed using ‘scaler’ and passed to ‘model’ for prediction. The train function reads in project_files.h5 which was created in the previous step and trains a logistic regression model using make_pipeline.

Moreover, a Python script CLASSIFIER_VISUALS.py was written for training and testing. The accelerometer data from project_files.h5. is read and loaded into data frames ‘Train’ and ‘Test’. Accelerometer data is separated and stored in X_train and X_test, which are normalized using

StandardScaler. The logistic regression model is set with a maximum iteration of 1000 and the pipeline is fitted to the training data. In this program the accuracy and recall scores are calculated and printed which is shown in Figure 28. A ROC curve is created using the predicted probabilities and the actual labels from test data. The area under the ROC curve, the AUC is calculated and printed which is also visible in Figure 28. Finally, a confusion matrix is developed as well, which is shown in Figure 30.

```
y_pred is: [0. 1. 0. 0. 1. 0. 1. 0. 0. 0. 1. 1. 0. 0. 1. 1. 1.]
y_clf_prob is: [[9.0751208e-01 9.2487996e-02]
[1.60073895e-03 9.9839261e-01]
[9.91423473e-01 8.57652726e-03]
[9.73944546e-01 2.60554542e-02]
[2.82729535e-02 9.7177047e-01]
[8.89742738e-01 1.10257262e-01]
[1.81965595e-04 9.99818034e-01]
[9.05957529e-04 9.9994042e-01]
[7.52396148e-01 2.47683852e-01]
[5.70366207e-01 4.29633793e-01]
[9.91391093e-01 8.60890652e-03]
[8.96903268e-01 1.03696792e-01]
[3.61818873e-03 9.96381811e-01]
[3.25641972e-02 9.67435803e-01]
[9.84493277e-01 1.55867230e-02]
[9.89166862e-01 1.08331375e-02]
[3.35967857e-03 9.96646321e-01]
[4.91731529e-01 5.08268471e-01]
[2.14853742e-03 9.97851465e-01]]
Accuracy is: 0.9473684210526315
0.9
AUC is: 1.0
```

Figure 28: Output of CLASSIFIER_VISUALS.py displaying accuracy and AUC.

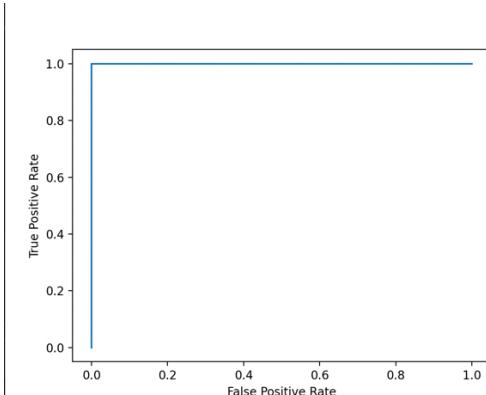


Figure 29: ROC developed from CLASSIFIER_VISUALS.py.

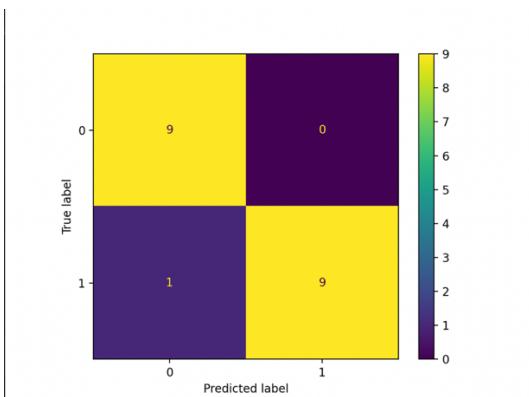


Figure 30: Confusion matrix developed by CLASSIFIER_VISUALS.py.

The code to accomplish these data and figures is displayed in Figure 31 and Figure 32.

```

1  import pandas as pd
2  import matplotlib.pyplot as plt
3  from sklearn.preprocessing import StandardScaler
4  from sklearn.model_selection import train_test_split
5  from sklearn.metrics import accuracy_score, recall_score, confusion_matrix, ConfusionMatrixDisplay, roc_curve, \
6  RocCurveDisplay, roc_auc_score
7  from sklearn.pipeline import make_pipeline
8  from sklearn.linear_model import LogisticRegression
9  from sklearn.inspection import DecisionBoundaryDisplay
10 from sklearn.decomposition import PCA
11 import h5py as h5
12 import numpy as np
13 import pandas as pd
14
15 # Step 2. Read the dataset as a Pandas data frame
16 with h5.File("./project_files.h5", "r") as hdf:
17     Train = pd.DataFrame(hdf.get('/dataset/train'))
18     Test = pd.DataFrame(hdf.get('/dataset/test'))
19
20 Train = Train.dropna()
21 Test = Test.dropna()
22
23 Train.columns = ["avgXAccel", "avgYAccel", "avgZAccel", "maxZAccel", "avgTimeAboveAvgZ", "avgAccel", "maxAccel", "avgTimeAboveAvg", "lat"]
24 Test.columns = ["avgXAccel", "avgYAccel", "avgZAccel", "maxYAccel", "avgZAccel", "maxZAccel", "avgTimeAboveAvgZ", "avgAccel", "maxAccel", "avgTimeAboveAvg", "lat"]
25
26 Y_train = Train['label']
27 Y_test = Test['label']
28
29 X_train = Train.iloc[:, :-1]
30 X_test = Test.iloc[:, :-1]
31
32
33 scaler = StandardScaler()

```

Figure 31: Python file CLASSIFIER_VISUALS.py part 1.

```

34 X_train = scaler.fit_transform(X_train)
35 X_test = scaler.transform(X_test)
36
37 # Now, similar to the code presented in the lecture slides, train and test your model.
38
39 lreg = LogisticRegression(max_iter=10000)
40 clf = make_pipeline(StandardScaler(), lreg)
41
42 clf.fit(X_train, Y_train)
43
44 y_pred = clf.predict(X_test)
45 y_clef_prob = clf.predict_proba(X_test)
46 print("y_pred is: ", y_pred)
47 print("y_clef_prob is: ", y_clef_prob)
48
49 acc = accuracy_score(Y_test, y_pred)
50 print("Accuracy is: ", acc)
51 recall = recall_score(Y_test, y_pred)
52 print(recall)
53
54 cm = confusion_matrix(Y_test, y_pred)
55 cm_display = ConfusionMatrixDisplay(cm).plot()
56 plt.show()
57
58 fpr, tpr, _ = roc_curve(Y_test, y_clef_prob[:, 1], pos_label=clf.classes_[1])
59 roc_display = RocCurveDisplay(fpr=fpr, tpr=tpr).plot()
60 plt.show()
61
62 auc = roc_auc_score(Y_test, y_clef_prob[:, 1])
63 print("AUC is: ", auc)
64

```

Figure 32: Python file CLASSIFIER_VISUALS.py part 2.

Model Deployment

The final step of the design process was to deploy the model in a desktop app. The goal was to have the application accept an input file in CSV format and generate a CSV file as the output, which would include the prediction label, either walking or jumping, for the selected input file. To build the graphical user interface the Python library Tkinter was used. The overall design is straight forward. The interface provides a select file button to choose a CSV file which would contain accelerometer data. The GUI also has buttons to then generate output based on data provided in the selected CSV. The Application's initializer method sets up the various buttons, frames, as well as a graph. It also loads the trained model from CLASSIFIER_FINAL.py. When the 'Select CSV File' button is clicked the user will see a file dialog box. When the file is chosen, it gets read to a DataFrame, and split into segments of 500 rows each. Both buttons to generate output are enabled when needed. The main window is shown in Figure 33.



Figure 33: Main window of simple GUI.

Upon generating the output, the prediction label and output graph are displayed. The prediction label will simply tell a user if the file they selected contained walking or jumping data. The figure generated shows a bar graph with each of the extracted features. An example of the generated output is shown below in Figure when a file called ‘jumping_backpocketL1.csv’ is uploaded to the interface. Though it should be noted that several tests were conducted with a wide variety of input files uploaded to ensure functionality. Another measure taken was to test walking and jumping data from individuals outside of the team. Meaning the model was successful for people who may have a different stride, higher jump etc. than the team members who built the model with their walking/jumping data.

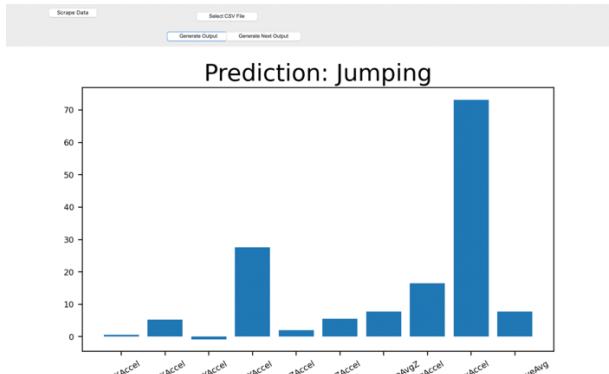


Figure 34: Main window after pressing the ‘Generate Output’ button with a file uploaded called ‘jumping_backpocketL1.csv’.

Real Time

After development of the model deployment, a real time implementation was attempted. The goal of this change would be to create an application which is capable of reading accelerometer data from a smart phone and classifying it in real time. To implement this a different button was created on the main window which allowed a user to scrape data from a web page.

```

1  import tkinter as tk
2  from tkinter import ttk as ttk
3  import pandas as pd
4  import matplotlib.pyplot as plt
5  from tkinter.filedialog import askopenfilename
6  from matplotlib.backends.backend_tkagg import FigureCanvasTkAgg
7  import classifier
8  import numpy as np
9  import math
10 import time
11 from bs4 import BeautifulSoup
12 from selenium import webdriver
13
14 class Application:
15
16     def __init__(self, master):
17         self.master = master
18         self.master.title("Walking/Jumping Classifier")
19
20
21         #scrape_button
22         self.scrape_button = ttk.Button(master, text="Scrape Data", command=self.scrape_data)
23         self.scrape_button.place(x=100, y=5)

```

Figure 35: Initialization of necessary libraries and a new button on the main window.

The approach taken was to enable remote access through the Phyphox app which is used for the data collection process. The smart phone on which data is being collected can then connect to a web page correlating to the IP address of that device this is seen in line 70 of Figure 36 as ‘<http://192.168.0.212>’. In order to scrape information from the Phyphox page the Python libraries Selenium and Beautiful Soup are included. Upon clicking the scrape button, a method is called initialize a Selenium web driver retrieving the Phyphox page. The BeautifulSoup library is used to try collecting the necessary data which can be put into a data frame. A while loop continuously extracts the latest data from the page. Lastly, using the previously established classifier a prediction is returned to tell a user which motion is being demonstrated. With further time and iteration, the real time implementation could have been fully developed.

```

66     def scrape_data(self):
67         # Web driver
68         driver = webdriver.Chrome()
69         # link correlating to the IP of the smart-phone
70         driver.get("https://192.168.0.212")
71         # load delay
72         time.sleep(1)
73
74         # Finding the necessary data
75         soup = BeautifulSoup(driver.page_source, "html.parser")
76         table = soup.find("table", {"id": "data-table"})
77         data = []
78
79         for row in table.find_all("tr"):
80             cols = row.find_all("td")
81             cols = [col.text.strip() for col in cols]
82             data.append(cols)
83
84         while True:
85             # Extract the latest data
86             html = driver.page_source
87             soup = BeautifulSoup(html, "html.parser")
88             table = soup.find("table", {"id": "data-table"})
89             last_data = table.find_all("tr")[-1]
90             last_values = [float(col.text.strip()) for col in last_data.find_all("td")[1:]]
91
92             # Create a DataFrame df with the most recent data
93             df = pd.DataFrame([last_values], columns=["Time", "X-Accel", "Y-Accel", "Z-Accel"])
94             df = df.rolling(window=4).mean().dropna()

```

Figure 36: Part 1 of Python code implemented to attempt real time classification.

```

94
95         # Delay
96         time.sleep(0.1)
97         # Classify data
98         pred = classifier.Classify(self.model, self.scaler, df)
99
100
101         # Print the predicted motion
102         if pred[0] == 1:
103             print("Walking")
104         else:
105             print("Jumping")
106
107         print(pred)

```

Figure 37: Part 2 of Python code implemented to attempt real time classification.

Participation Report

A summary of group participation can be seen below in Table 1.

Table 1: Participation table summarizing contributions.

| Program | |
|-------------------------|----------------------|
| Section | Member(s) |
| Data Collection | Ethan, Lauren, Jacob |
| Data Storage | Ethan, Lauren, Jacob |
| Visualization | Ethan |
| Pre-processing | Ethan, Jacob |
| Feature Extraction | Ethan, Jacob |
| Creating a Classifier | Ethan, Jacob |
| Deployment | Lauren, Jacob |
| Real Time | Lauren |
| Report | |
| Section | Member(s) |
| Data Collection | Lauren |
| Data Storage | Lauren, Ethan |
| Visualization | Lauren, Ethan |
| Pre-processing | Lauren, Jacob, Ethan |
| Feature Extraction | Lauren, Jacob |
| Training the Classifier | Lauren, Jacob |
| Deployment | Lauren |
| Real Time | Lauren |
| Video | |
| Section | Member(s) |
| Recording | Ethan, Jacob, Lauren |
| Editing | Ethan |

Bibliography

- [1] A. Etemad, "ELEC 390 Principles of Design & Dev. W23 Slides - Part 1," 2023. [Online]. Available: <https://onq.queensu.ca/d2l/le/content/711190/viewContent/4584784/View>. [Accessed 7 April 2023].
- [2] A. Etemad, "ELEC 390 Principles of Design & Dev. W23 Slides - Part 2," 2023. [Online]. Available: <https://onq.queensu.ca/d2l/le/content/711190/viewContent/4618175/View>. [Accessed 7 April 2023].
- [3] A. Etemad, "ELEC 390 Principles of Design & Dev. W23 Project Instructions," 2023. [Online]. Available: <https://onq.queensu.ca/d2l/le/content/711190/viewContent/4586805/View>. [Accessed 24 Feb 2023].
- [4] A. Etemad, "ELEC 390 Principles of Design & Dev. W23 Slides - Part 3," 2023. [Online]. Available: <https://onq.queensu.ca/d2l/le/content/711190/viewContent/4634633/View>. [Accessed 7 April 2023].
- [5] A. Etemad, "ELEC 390 Principles of Design & Dev. W23 Slides - Part 4," 2023. [Online]. Available: <https://onq.queensu.ca/d2l/le/content/711190/viewContent/4647961/View>. [Accessed 7 April 2023].
- [6] p. Ippolito, "Feature Extraction Techniques," Medium, 10 October 2019. [Online]. Available: <https://towardsdatascience.com/feature-extraction-techniques-d619b56e31be>. [Accessed 12 April 2023].
- [7] P. Marcelino, "Data Analysis and Feature Extraction with Python," Kaggle, 2022. [Online]. Available: <https://www.kaggle.com/code/pmarcelino/data-analysis-and-feature-extraction-with-python>. [Accessed 12 April 2023].