

## Title: Task B-1 Report

**Author:** Robin Findlay-Marks, s103603871

### Setup Attempt:

Option B - Task 2: Data processing 1 took a long time to complete and had a lot of trouble understanding what I was meant to do.

The task seems to suggest that it is recommended for me to copy the code from P1 and then explain how it works. However P1 seems to do thing so differently from v0.1 that it is impossible without overhauling v0.1 to the point where it is basically P1 itself. I instead decided to figure out how to get the loading and processing done myself. I struggled for a while to get it to work properly but I'm fairly happy with what I did.

I made the program download a single dataframe, rather than two, which is defined with two dates. This dataframe is then split using one of three methods, with a third date, a ratio or simply randomly. I did this by having several smaller functions that first of all handle which splitting method to use, by checking which parameters are set, and second of all having functions dealing with the splitting by ratio or date.

The NaN issue was a annoying. I didn't understand what I meant by the "NaN issue in the data" as I couldn't find it referenced anywhere. I did see that something about NaNs was referenced in the P1 code, and after I realised that I researched what NaNs are and used the line from P1 to remove them.

Task 2 – 1.e was very confusing, and I'm not sure if it's implemented or not. It seems to want me to scale the feature columns and store them as data structure, but it already scales the feature columns and stores them as an array, which is a data structure. So as far as I can tell, this is already done, but that doesn't seem right.

The rest of the explanations of how my code works can be found as comments in the code.

I have included screenshots for easy access.

### Parameters:

```
44 #
45 #
46 #
47 # -----Parameters-----
48 #
49 #
50 #
51 #
52 DATA_SOURCE = "yahoo"
53 COMPANY = "TSLA"
54
55 TRAIN_START = '2015-01-01' #Start date of dataset #Must be in 'YYYY-MM-DD' format eg '2015-01-01'
56 TEST_END = '2022-12-31' #End date of dataset #Must be in 'YYYY-MM-DD' format eg '2022-12-31'
57
58 SPLIT_DATE = '2020-01-01' #Split date of dataset #Must be in 'YYYY-MM-DD' format eg '2020-01-01'
59 SPLIT_DATE_BOOL = False
60
61 RATIO = 4 #Int or Float #Not actually a ration, but idk what else to call it
62 RATIO_BOOL = False #2 is train/test equally split, 4 is train gets about 75% of data
63
64 #If both SPLIT_DATE_BOOL and RATIO_BOOL are false, it picks a random date
65
66 SCALER = True # Pick whether to scale feature columns or not
67
68 # Train and test data global variables for setting
69 trainData = None
70 testData = None
```

## Downloading and storing of data as file:

```

72 # function for checking if the data is already downloaded (needs internet connection to download)
73 # if the data is NOT in a file, it downloads the data and makes a csv file and returns the data
74 # if the data IS in a file, it reads the data from the file and returns the data
75 def checkFiles(filename):
76     if (os.path.exists(filename)):
77         #Read csv file and return the data inside
78         data = pd.read_csv(filename)
79         # NaN values from pandas are values that are not present. For example in stocks if the stock data for a
80         # specific day was not recorded, I believe it would still have a record for that day, only the values
81         # would be NaN or 'Not a Number'
82
83         #what dropna() does is simply remove the missing values from the dataset
84         data.dropna(inplace=True)
85         return data
86     else:
87         #download data from online
88         data = yf.download(COMPANY, start=TRAIN_START, end=TEST_END, progress=False)
89
90         # Save data to csv file
91         data.to_csv(filename)
92         # For some reason it needs to read it from the file otherwise it won't work
93         data = pd.read_csv(filename)
94         # remove NaN values from the dataset
95         data.dropna(inplace=True)
96         return data

```

## Splitting data by date

```

119 # this function gets the datafile name as well as the split date
120 # it then runs the file checker to get the dataset, then splits the dataset at the split date
121 def getDataSplitDate(filename, splitDate):
122     df = checkFiles(filename)
123
124     # Make it know that the date column is indeed a date
125     df['Date'] = pd.to_datetime(df['Date'])
126     df = df.set_index(df['Date'])
127     df = df.sort_index()
128
129     #Convert input to datetime, add 1 day, then convert back to string+
130     date = datetime.strptime(splitDate, '%Y-%m-%d')
131     testStartDate = date + timedelta(days=1)
132     testStartDate = testStartDate.strftime('%Y-%m-%d')
133
134     # create train/test partition
135     global trainData
136     trainData = df[TRAIN_START:splitDate]
137     global testData
138     testData = df[testStartDate:TEST_END]
139     print('Train Dataset:', trainData.shape)
140     print('Test Dataset:', testData.shape)
141
142     #trainData.to_csv("trainfilename.csv")
143     #testData.to_csv("testfilename.csv") #test that the data is split correctly
144

```

## Splitting data by ratio

```

# this function gets the datafile name as well as the ratio number
# it then runs the file checker to get the dataset, then splits the dataset at the split date
def getDataRatio(filename, ratio):
    df = checkFiles(filename)

    # Make it know that the date column is indeed a date
    df['Date'] = pd.to_datetime(df['Date'])
    df = df.set_index(df['Date'])
    df = df.sort_index()

    # Convert strings to dates
    date1 = datetime.strptime(TRAIN_START, '%Y-%m-%d')
    date2 = datetime.strptime(TEST_END, '%Y-%m-%d')

    # do math to get the date we want
    trainEndDate = date2 + (date1 - date2) / ratio

    #Convert input to datetime, add 1 day, then convert back to string+
    print("Middle : " + trainEndDate.strftime('%Y-%m-%d'))
    testStartDate = trainEndDate + timedelta(days=1)
    #testStartDate = testStartDate.strftime('%Y-%m-%d') # i don't remember why i commented this, but it works

    # create train/test partition
    global trainData
    trainData = df[TRAIN_START:trainEndDate]
    global testData
    testData = df[testStartDate:TEST_END]

    #trainData.to_csv("trainfilename.csv")
    #testData.to_csv("testfilename.csv") #test that the data is split correctly

```

## Processing of parameters

```

177 def getData(): #Main function for deciding which split method was chosen
178     #Make filename for the saved data file
179     ticker_data_filename = os.path.join("data", f"{COMPANY}_{TRAIN_START}_{TEST_END}.csv")
180     if (SPLIT_DATE_BOOL):
181         getDataSplitDate(ticker_data_filename, SPLIT_DATE)
182     elif (RATIO_BOOL):
183         getDataRatio(ticker_data_filename, 3.5)
184     else: #Random Date
185         #Convert dates to datetime
186         dateStart = datetime.strptime(TRAIN_START, '%Y-%m-%d')
187         dateEnd = datetime.strptime(TEST_END, '%Y-%m-%d')
188         #Get random date inbetween the start and end
189         random_date = dateStart + (dateEnd - dateStart) * random.random()
190         #convert back to string
191         random_date = random_date.strftime('%Y-%m-%d')
192         #Use random date as split
193         getDataSplitDate(ticker_data_filename, random_date)
194
195     getData()

```