# Assignment 1

DATA MINING

CSE 572: Spring 2020

**Submitted by**
**Akshay Kumar(akuma216@asu.edu)**

# 1. Assignment Phase 1: Feature Extraction

As a part of assignment 1 we are given CGMData and CGMTime file to extract 4 features that are useful for the prediction process. I have started by looking into 5 different types of feature extraction techniques that can be useful for this time series data.

1. Fast Fourier Transform
2. Discrete Wavelet Transform
3. Velocity
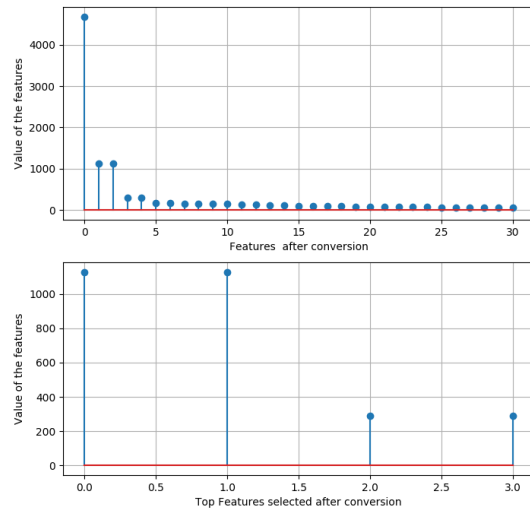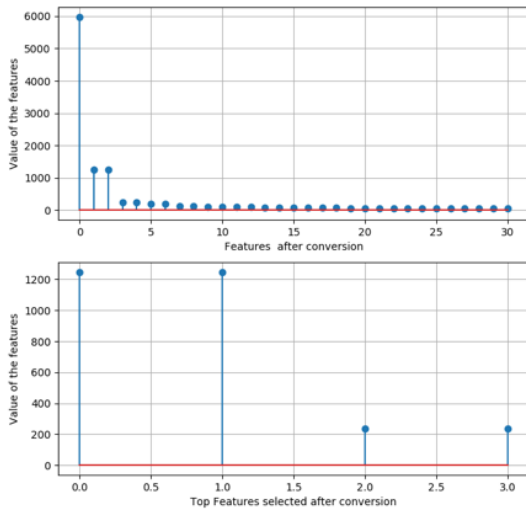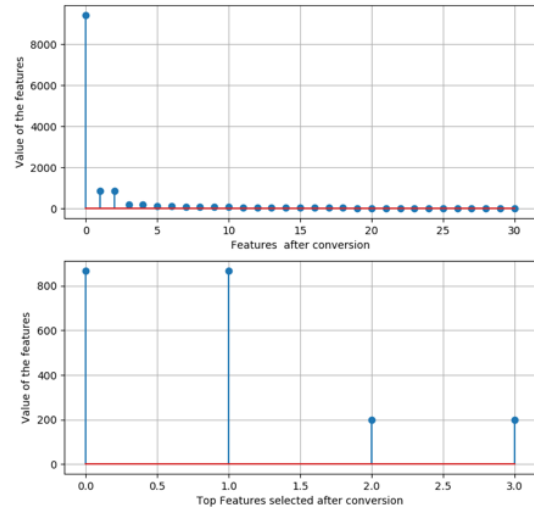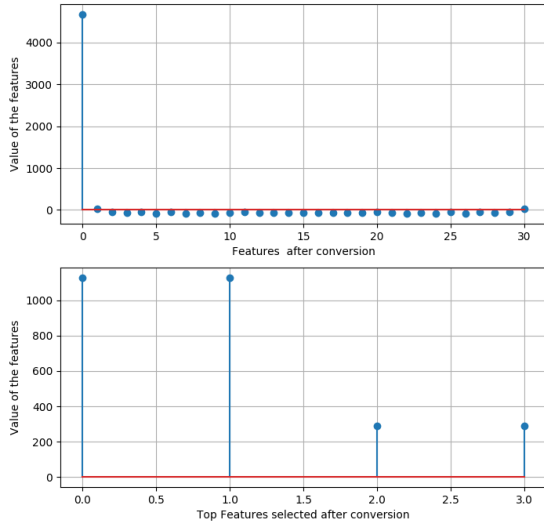4. Moving Average
5. Expanding Window

## Question 1 and 2
## 1.FAST FOURIER TRANSFORM

Fast Fourier Transform is used to extract frequency-based information from the time series data. We only select those frequency components which show high variance and keep it as a good feature among the whole array. For this project I have kept top 4 features for each row of data. I also have ignored the first highest feature as it comes from a constant 1 which act as a noise. I have plotted the original extracted data in **decreasing order** with the selected **top 4 data** on the bottom.

**Reason for choosing FFT**
It gives the peak of the CGM data at which meal was taken showing the rise of glucose level in the body. For now, this can be considered as a good feature for detecting the meal intake in the graph

All 31 graph for each and every person can be found under Person folder and subfolder FFT.

## 2.DISCRETE WAVELET TRANSFORM

Discrete wavelet transform uses different sets of wavelet scales and transition to find the frequency domain of the given series. It is more helpful in finding the change as compared to Fourier transform because it also captures the rate and location. For this also I have taken top 4 features as my feature matrix

**Reason for choosing DWT**

Discrete Wavelet Transform is similar to FFT but the advantage of choosing DWT is that it capture the rate of change more accurately as compared to the FFT and also the time stamp if we want to know the exact location.



All 31 graph for each and every person can be found under Person folder and subfolder DWT.
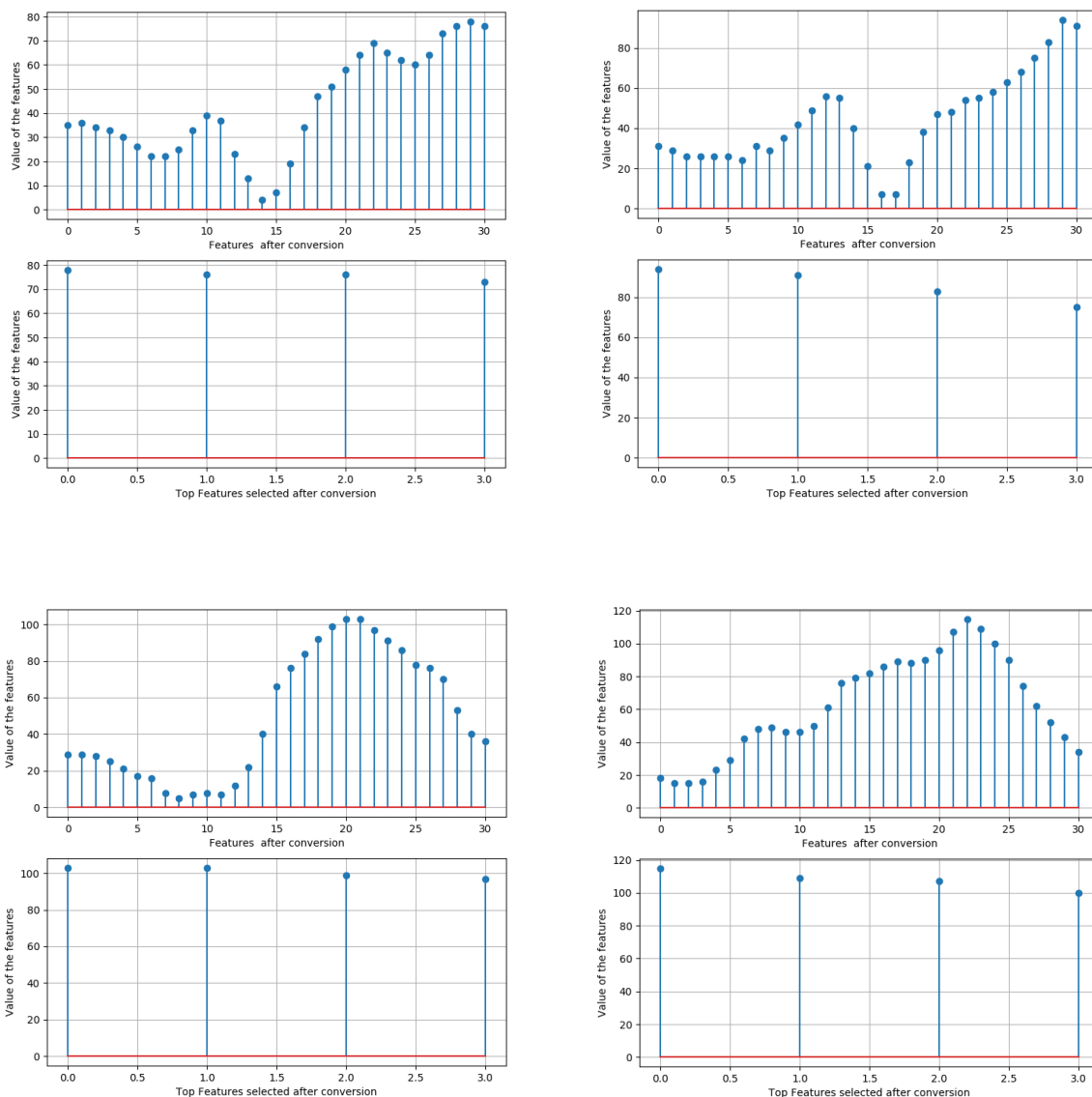
## 3.VELOCITY

Velocity finds the rate of change between the t+1 and t timestamp of the time series data. Its find the rate of change of the value between these points.
For this also I have taken top 4 features as my feature matrix

**Reason for choosing Velocity**

It captures rate of change of y value which will help us to locate the maximum point and points around it and I selected the top 4 of those points as my features



All 31 graph for each and every person can be found under Person folder and subfolder Velocity.

## 4. MOVING AVERAGE

Moving average is used to create smooth version of the time series data set. It helps to remove any random noise that can be present over the series or any random fluctuations. I have used rolling function of pandas creating moving average for length 16. I have removed the NAN values in order to create plots and finding the maximum 4 values. For this also I have taken top 4 features as my feature matrix

**Reason for choosing Moving Average**

Moving average smoothes the data and form a good features for creating a stationary points also keeping previous information.



All 31 graph for each and every person can be found under Person folder and subfolder RollingWindow.

# 5. EXPANDING WINDOW

Expanding windows is same as cumulative sum which keeps the record of all the past things and increases over time. Same happens in expanding window where number of observation increases over time.

All 31 graph for each and every person can be found under Person folder and subfolder ExpandingWindow.

# QUESTION 3

Drawbacks and Advantages of choosing all the 5 features above

1. Fast Fourier Transform- We have used FFT to capture the increase in the glucose level of a person when there is food intake. According to initial hypothesis it should have selected the points with high change in value but in some of the case it selected peak from graph with multiple maxima. It

doesn't capture the exact point where there is a spike but can also capture multiple points where there is not sudden increase. The other thing is we have to deal with real and imaginary points while dealing with FFT. Since most of the dataset has one single maxima it forms a decent feature. To avoid this I used DWT to look further

2. Discrete Wavelet Transform- In DWT both higher frequency window and lower frequency window helps to capture rate of change in glucose level more accurately as compared to the FFT.

3. Velocity- It tries to find the maximum point in the graph but sometimes lead to find the maximum point where these is least rate of change in the glucose level.

4. Moving Average- Its tries to smooth the graph and create stationary points with minimum change in the glucose level. It provides a constant value keeping the value close to the maxima if window size is selected properly.

5. Expanding Window- Expanding window is not at all helpful in this case as it keeps all the information from the past and doesn't help in finding the exact instantaneous change in the value.

# Question 4
# FEATURE MATRIX

I have created a FeatureExtracted.csv file which contains all the features that I have extracted using the above feature extraction technique. Since I have extracted 5 features having top 4 from each method we have matrix of size 32X20 for each person
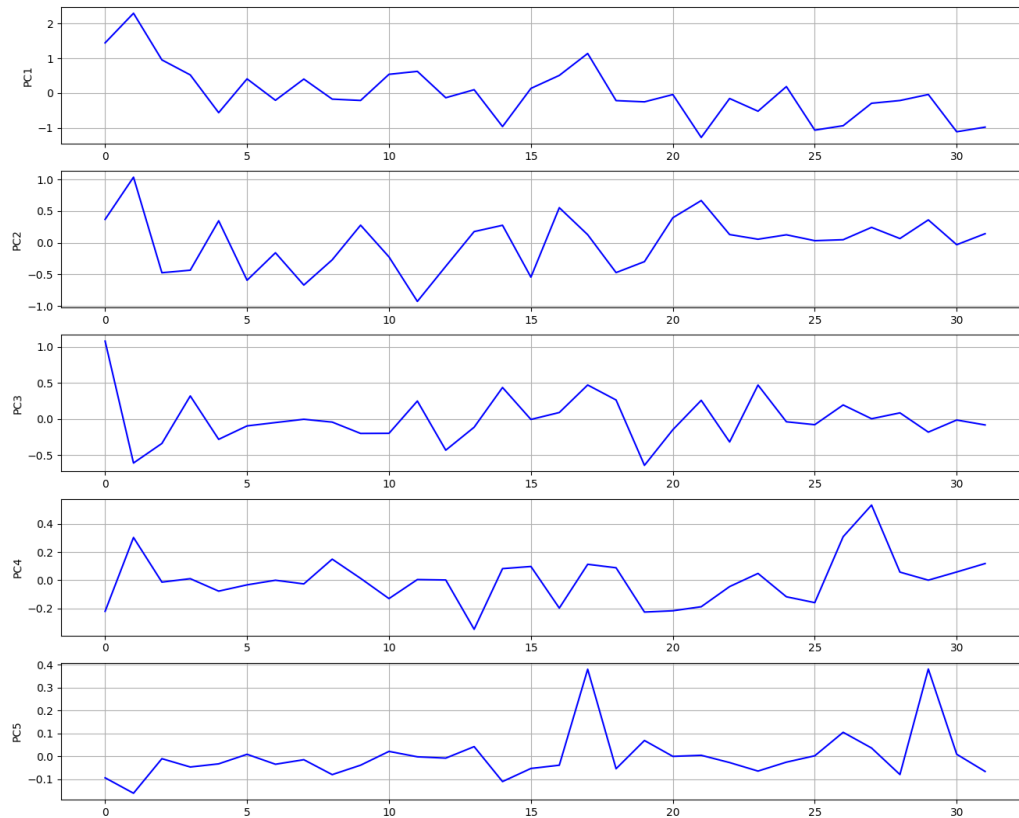
Screenshot of the feature matrix

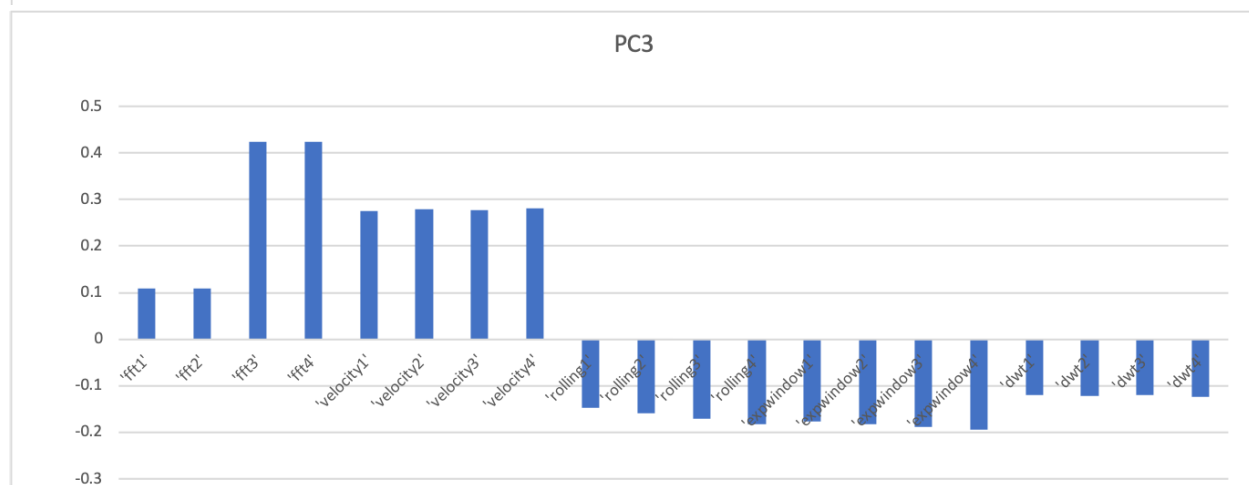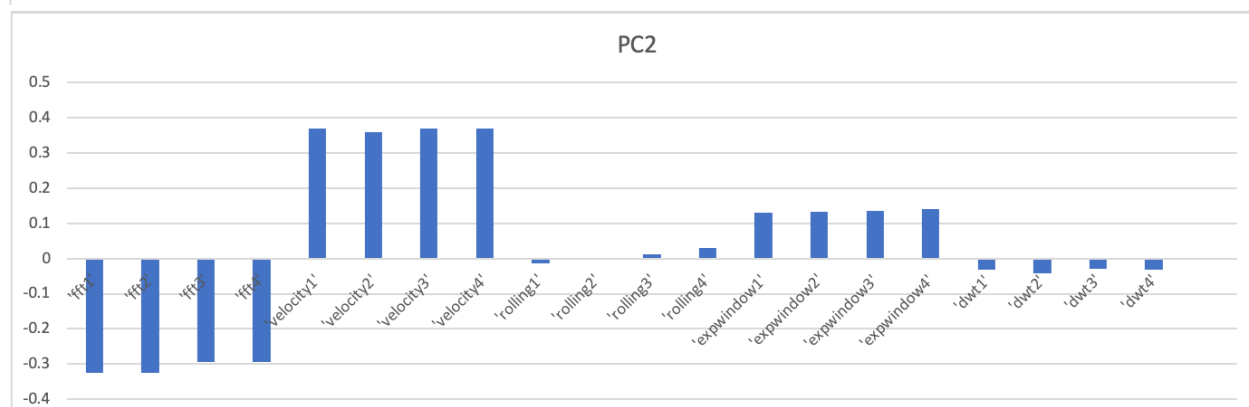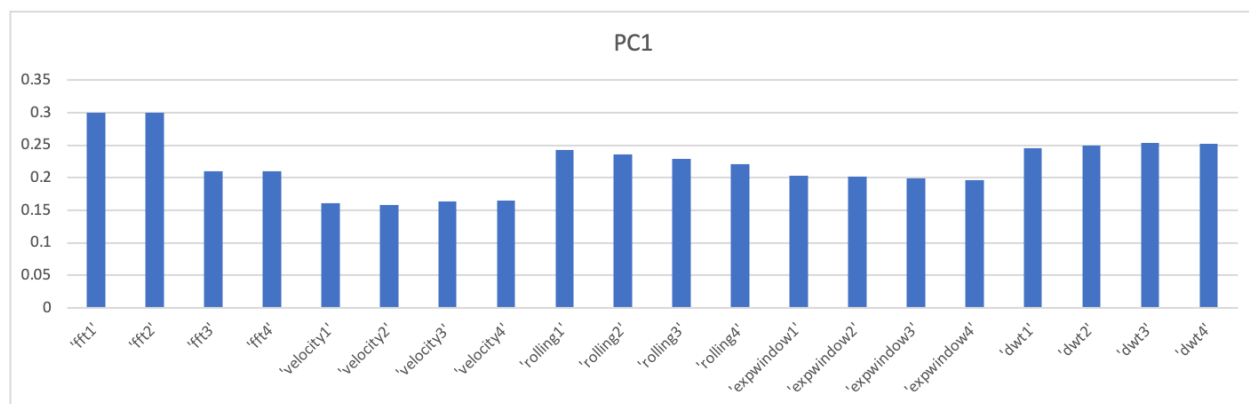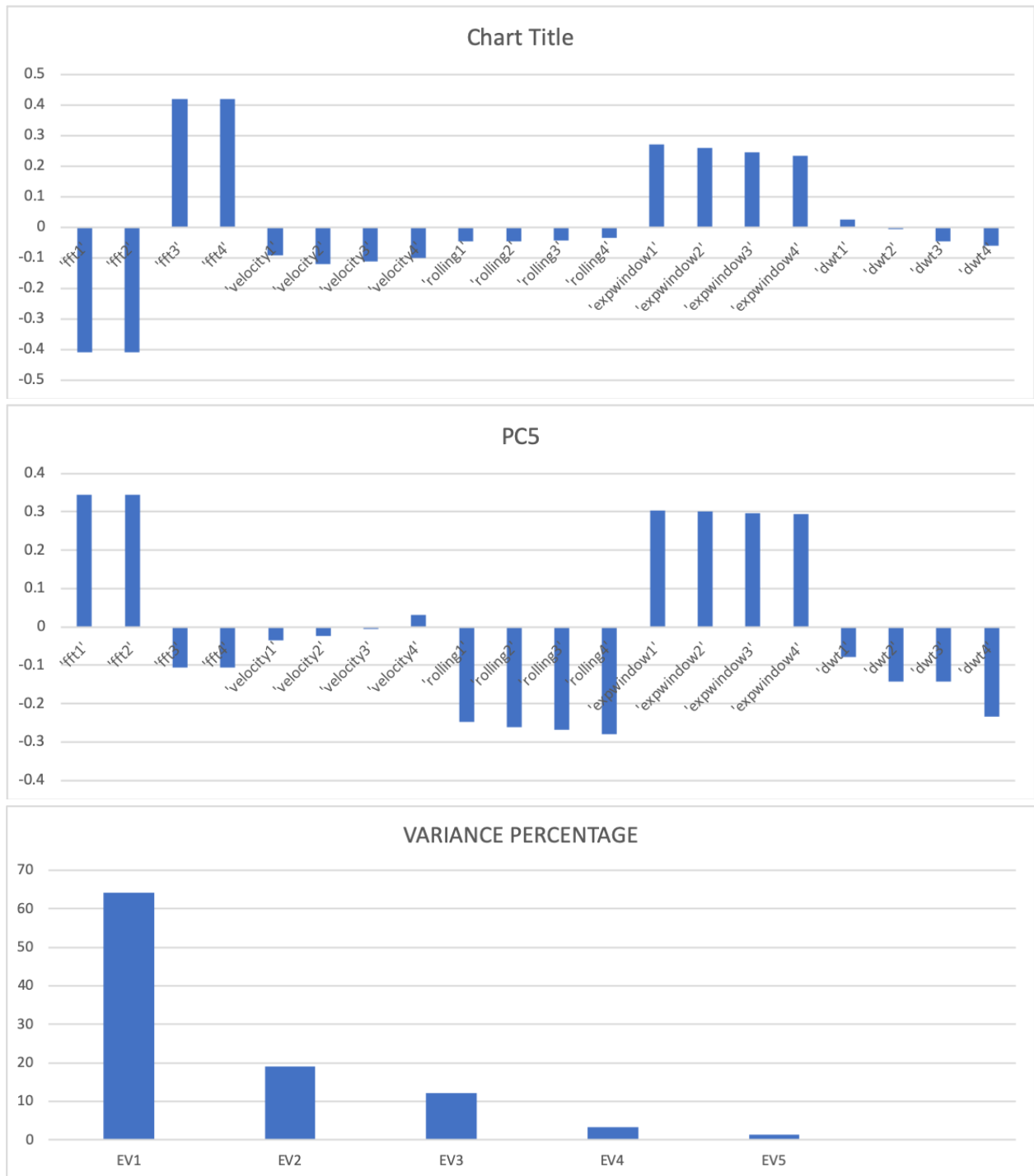| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1349.4563506679900 | 1349.4563506679900 | 344.274083297032000 | 344.274083297032000 | 233.0 | 232.0 | 225.0 | 217.0 | 229.125 | 222.25 | 214.3125 | 205.125 | 167.6129032258060000 | 164.666666666670000 | 161.4482758620690000 | 157.9285714128571000 | 366.281312654630000 | 362.038671967512000 | 362.038671967512000 | 352.139177030900100 |
| 1 | 868.543990498082000 | 868.543990498082000 | 199.260485819247000 | 199.260485819247000 | 205.0 | 198.0 | 198.0 | 187.0 | 338.125 | 338.125 | 336.8125 | 336.75 | 308.153846153846000 | 308.08 | 307.962962962963 | 307.75 | 494.974746830583000 | 492.146319705837000 | 490.024999362277000 | 489.317892581091000 |
| 2 | 1247.374320481750000 | 1247.374320481750000 | 235.889208703414000 | 235.889208703414000 | 61.0 | 60.0 | 59.0 | 58.0 | 248.125 | 244.625 | 239.1875 | 232.4375 | 192.935483870968000 | 191.633333333330000 | 190.069655172410000 | 188.0 | 383.958982184295000 | 381.837661840736000 | 370.523953341751000 | 367.695526217005000 |
| 3 | 1128.069701321870000 | 1128.069701321870000 | 290.575521861880000 | 290.575521861880000 | 94.0 | 91.0 | 83.0 | 75.0 | 202.375 | 195.8125 | 187.75 | 179.8125 | 150.935483870968000 | 148.366666666670000 | 145.517241379310000 | 142.785714285714000 | 322.440692210680 | 320.319371875060000 | 304.763022691402 | 302.641702347842000 |
| 4 | 359.871193700071000 | 359.871193700071000 | 74.397787707423600 | 74.397787707423600 | 78.0 | 76.0 | 76.0 | 73.0 | 157.625 | 157.5 | 156.9375 | 156.75 | 143.481481481482000 | 143.392857142857000 | 143.346153846154000 | 143.241379310340500 | 234.759451353930400 | 232.638131010374000 | 229.102597104441000 | 224.859956417322000 |
| 5 | 1147.917950242590000 | 1147.917950242590000 | 240.674122121889000 | 240.674122121889000 | 49.0 | 48.0 | 45.0 | 44.0 | 210.3125 | 206.9375 | 202.5 | 197.0625 | 158.677419354839000 | 156.866666666670000 | 154.862068965517000 | 152.714285714286000 | 316.783837971573000 | 311.126983722081000 | 304.055915910215000 | 303.348800912902900 |
| 6 | 697.734983029905000 | 697.734983029905000 | 176.909565549206000 | 176.909565549206000 | 64.0 | 63.0 | 60.0 | 57.0 | 172.125 | 170.125 | 166.3125 | 161.8125 | 140.258064516129000 | 139.466666666670000 | 138.517241379310000 | 137.392857142857000 | 263.043722601396000 | 260.215295476650000 | 255.265548008344000 | 243.244732728172000 |
| 7 | 1162.919247676770000 | 1162.919247676770000 | 258.031447638800010 | 258.031447638800010 | 47.0 | 46.0 | 45.0 | 45.0 | 203.5 | 196.875 | 189.1875 | 180.9375 | 151.741935483871000 | 149.2 | 146.413793103484000 | 143.428571428571000 | 325.269119345812000 | 322.440692210660 | 321.733585439879000 | 319.612650963200000 |
| 8 | 660.112104154348000 | 660.112104154348000 | 208.879607155198000 | 208.879607155198000 | 53.0 | 53.0 | 50.0 | 47.0 | 175.5 | 172.625 | 168.9375 | 164.75 | 144.451612903226000 | 143.033333333330000 | 141.551724137931000 | 140.0 | 264.457936163769000 | 266.579256507328000 | 259.508188695463 | 253.851334446597100 |
| 9 | 480.573716064795000 | 480.573716064795000 | 122.301190056877000 | 122.301190056877000 | 86.0 | 85.0 | 85.0 | 84.0 | 171.25 | 171.0625 | 170.5 | 169.8125 | 157.086956521739000 | 157.045454545454000 | 156.833333333333000 | 156.4 | 272.236110756821000 | 266.579256507328000 | 262.336615820200900 | 251.022907321224000 |
| 10 | 1088.284162532940000 | 1088.284162532940000 | 192.437697229543000 | 192.437697229543000 | 85.0 | 79.0 | 71.0 | 66.0 | 221.25 | 218.9375 | 215.1875 | 210.25 | 173.935483870968000 | 172.433333333333000 | 170.827586206897000 | 169.071428571429000 | 323.854065783439000 | 321.733585439879000 | 321.733585439879000 | 314.662517628014000 |
| 11 | 1379.123632577880 | 1379.123632577880 | 327.576422117920 | 327.576422117920 | 47.0 | 47.0 | 47.0 | 47.0 | 205.9375 | 198.1875 | 189.5625 | 180.3125 | 146.387096774194000 | 143.7 | 145.8279661294300 | 153.512741379310000 | 151.035714285714000 | 282.135605693432000 | 280.721392131059000 | 276.478751442940 | 275.064537881567 |
| 12 | 773.131627455207 | 773.131627455207 | 152.111791197780000 | 152.111791197780000 | 30.0 | 29.0 | 21.0 | 20.0 | 189.25 | 186.9375 | 183.875 | 179.8125 | 154.870967741936000 | 153.8 | 152.517241379310000 | 151.035714285714000 | 282.135605693432000 | 280.721392131059000 | 276.478751442940 | 275.064537881567 |
| 13 | 879.758177030342000 | 879.758177030342000 | 114.089729406830000 | 114.089729406830000 | 103.0 | 103.0 | 99.0 | 97.0 | 187.8125 | 186.9375 | 184.625 | 181.0625 | 150.354838709677000 | 150.133333333333000 | 149.724137931034000 | 148.892857142857000 | 299.913780286486500 | 285.671139599365000 | 283.549819255806000 | 272.236110756821000 |
| 14 | 218.416995982810000 | 218.416995982810000 | 165.916947160289000 | 165.916947160289000 | 98.0 | 97.0 | 93.0 | 92.0 | 109.6875 | 108.5625 | 107.375 | 106.375 | 101.193548387097000 | 100.533333333330000 | 99.7931034482759000 | 99.0714285714286000 | 171.119841047145000 | 171.119841047145000 | 167.584307141212000 | 167.584307141212000 |
| 15 | 926.897721711028000 | 926.897721711028000 | 246.997845376230000 | 246.997845376230000 | 46.0 | 43.0 | 43.0 | 42.0 | 187.75 | 182.375 | 176.125 | 169.4375 | 147.354838709677000 | 145.366666666670000 | 143.103448275862000 | 140.607142857143000 | 305.470129472589000 | 299.813275223090600 | 295.570634535977000 | 292.742207411231000 |
| 16 | 795.096925420657000 | 795.096925420657000 | 152.796018088498000 | 152.796018088498000 | 153.0 | 152.0 | 151.0 | 150.0 | 210.4375 | 207.375 | 203.3125 | 198.3375 | 175.064516129032000 | 173.7 | 172.172413793103000 | 175.035714285714000 | 309.005663378521000 | 308.298556697335000 | 308.298556697335000 | 306.177236253775000 |
| 17 | 1299.737648696850000 | 1299.737648696850000 | 300.138460038501000 | 300.138460038501000 | 158.0 | 156.0 | 153.0 | 151.0 | 206.25 | 203.25 | 199.625 | 194.875 | 216.8 | 216.727272727273000 | 216.111111111111000 | 215.75 | 326.683329081085 | 318.905158315133 | 316.076731190087000 | 297.691954879537000 |
| 18 | 811.854060298789000 | 811.854060298789000 | 252.334462839519000 | 252.334462839519000 | 57.0 | 54.0 | 51.0 | 50.0 | 161.4375 | 157.4375 | 152.1875 | 146.0 | 123.645161290323000 | 121.766666666667000 | 119.689655172414000 | 117.5 | 256.679761570717000 | 254.558441227157000 | 244.659946290545000 | 226.981276760882000 |
| 19 | 817.658085759332000 | 817.658085759332000 | 82.556906192392700 | 82.556906192392700 | 23.0 | 20.0 | 18.0 | 17.0 | 183.25 | 182.5625 | 182.25 | 180.625 | 150.655172413793000 | 150.607142857143000 | 150.433333333333000 | 150.259259259259000 | 287.792459429250 | 285.671139599365000 | 275.064537881567000 | 266.579256507328000 |
| 20 | 644.632957401331000 | 644.632957401331000 | 99.987216964085500 | 99.987216964085500 | 115.0 | 109.0 | 107.0 | 100.0 | 180.75 | 180.6875 | 179.3125 | 178.75 | 155.928571428571000 | 150.925925925926000 | 155.461538461538000 | 285.671139599365000 | 270.821897194448000 | 265.165042944955000 | 242.537629449800 | |
| 21 | 103.620392103020000 | 103.620392103020000 | 55.469757433287800 | 55.469757433287800 | 113.0 | 111.0 | 110.0 | 105.0 | 93.625 | 93.5 | 93.3125 | 93.0 | 100.75 | 99.0 | 96.833333333333000 | 94.714285714285700 | 144.249783620560 | 140.714249456123000 | 140.714249456123000 | 140.007142674936000 |
| 22 | 582.745653053067000 | 582.745653053067000 | 116.130652900043000 | 116.130652900043000 | 79.0 | 74.0 | 64.0 | 64.0 | 179.5625 | 179.0 | 177.25 | 174.75 | 157.433333333333000 | 157.379310344828000 | 157.193548387097000 | 157.035714285714000 | 284.964032818179000 | 283.549819255806000 | 271.529003975634000 | 261.629509039023000 |
| 23 | 528.518407223770000 | 528.518407223770000 | 210.688838013881000 | 210.688838013881000 | 105.0 | 101.0 | 91.0 | 86.0 | 134.1875 | 133.5625 | 132.75 | 131.0 | 113.107142857143000 | 113.037037037037000 | 112.758620689655000 | 112.423076923077000 | 204.353859762912000 | 201.525432638166000 | 197.989898732233000 | 194.454364826301000 |
| 24 | 796.153756962939000 | 796.153756962939000 | 163.280009226565000 | 163.280009226565000 | 96.0 | 96.0 | 96.0 | 96.0 | 193.5 | 191.9375 | 188.75 | 184.5 | 158.451612903226000 | 157.866666666670000 | 156.896551724138000 | 156.821428571429000 | 299.813275223090 | 287.792459429250 | 287.085353161780000 | 280.014285349873000 |
| 25 | 385.766567298227000 | 385.766567298227000 | 84.345331797752100 | 84.345331797752100 | 60.0 | 58.0 | 45.0 | 32.0 | 118.5625 | 118.25 | 117.9375 | 117.1875 | 103.574128571429000 | 103.518518518519000 | 103.517241379310000 | 103.384615384615000 | 188.797510576808000 | 184.554869886989000 | 177.483802078230000 | 171.826947828331000 |
| 26 | 261.412790613670000 | 261.412790613670000 | 176.871238487650 | 176.871238487650 | 61.0 | 59.0 | 59.0 | 58.0 | 100.125 | 99.875 | 99.875 | 99.4375 | 136.5 | 131.2 | 126.0 | 121.285714285714000 | 209.303607231218000 | 176.776692596370 | 161.927452897190 | 161.220346110533000 |
| 27 | 255.839781655207000 | 255.839781655207000 | 205.599103813674000 | 205.599103813674000 | 92.0 | 80.0 | 74.0 | 74.0 | 152.4375 | 149.375 | 146.6875 | 144.5 | 181.25 | 178.6 | 175.333333333333000 | 171.714285714286000 | 261.629509039023000 | 251.022907321224000 | 231.223917448001000 | 210.010714012405000 |
| 28 | 568.440767881526000 | 568.440767881526000 | 181.752850314689000 | 181.752850314689000 | 87.0 | 87.0 | 84.0 | 80.0 | 168.6875 | 164.75 | 160.9375 | 157.3125 | 141.838709677419000 | 140.233333333330000 | 138.793103448276000 | 137.464285714286000 | 268.700576850888000 | 253.144227664784000 | 242.537629449860000 | 241.830519165799000 |
| 29 | 698.225740477790000 | 698.225740477790000 | 112.319477948003000 | 112.319477948003000 | 98.0 | 98.0 | 94.0 | 92.0 | 157.625 | 151.0625 | 144.625 | 138.75 | 192.75 | 191.4 | 190.666666666667000 | 189.428571428571000 | 265.165042944955000 | 242.537629449800 | | |
| 30 | 290.523466950697000 | 290.523466950697000 | 121.097197760756000 | 121.097197760756000 | 43.0 | 42.0 | 41.0 | 38.0 | 110.3125 | 110.25 | 110.0625 | 109.6875 | 109.0 | 108.933333333330000 | 108.764705882353000 | 108.428571428571000 | 185.969083452063000 | 175.362481734264000 | 170.412734265958000 | 159.806132548160000 |
| 31 | 200.064783104811000 | 200.064783104811000 | 116.192022917808000 | 116.192022917808000 | 58.0 | 53.0 | 52.0 | 51.0 | 127.375 | 127.375 | 127.375 | 127.375 | 121.214285714286000 | 121.185185185185000 | 121.0 | 120.862068965517000 | 188.090403795622000 | 185.261976607656000 | 181.019335983756000 | 181.019335983756000 |

# Question 5
# PRINCIPAL COMPONENT ANALYSIS

For principal component analysis I am using the extracted feature matrix as input for PCA function from sklearn. It outputs an array of 31X5 matrix consisting of 5 principal components. Along with that I have created an eigen value matrix of 5X20 to see the contribution of each extracted feature on the 5 components Plots of eigen value matrix. Plotted all the five components in time series below.

PC1



PC2



PC3

## Question 6

First principal component shows FFT, Moving Average and DFT as the most important feature. The variance for PC1 is around 63% among all the principal components

The second principal component shows Velocity and Expanding window as the most important feature. The variance for PC1 is around 19% among all the principal components

The third principal component shows Velocity and FFT as the most important feature. The variance for PC3 is around 11% among all the principal components

The fourth principal component shows Expanding window as the most important feature. The variance for PC4 is around 3% among all the principal components

Similarly it can be seen for all the Principal components in the figure plotted above. I have also made a heatmap using seaborn library for analysis.