

PRINCIPAL COMPONENT ANALYSIS

Vishal S, IIT Madras

1 Introduction

In this assignment we will try to implement Principal Component Analysis method to a data and try to analyze the variance of the data with respect to its principal components

2 Principal Component Analysis

Principal Component Analysis is a procedure by which the data is reduced to a set of variables which are linearly uncorrelated called as **Principal Components**. These principal components are orthogonal to each other. These components are so chosen such that the variance is maximum in that direction.

3 Code

A code was written in MATLAB to calculate the principal components of a given data.

```
1 mean = zeros(len(2),1);
2 sum_col = zeros(len(2),1);
3 % To make the mean zero
4 for j = 1:len(2)
5     sum_col(j) = sum(fatiguedatademo(:,j));
6     mean(j) = sum_col(j)/len(1);
7 end
8 for j = 1:len(2)
9     for i = 1:len(1)
10        fatiguedatademo(i,j) = fatiguedatademo(i,j) - mean(j);
11    end
12 end
13
14 %% Calculating the principal components
15 [U,S,V] = svd(fatiguedatademo','econ');
16 error = zeros(26,1);
17 %% Calculating the error while reconstructing data
18 for k = 1:26
19     num_comp = k;
20     c1 = U(:,1:num_comp)*fatiguedatademo';
21     temp = c1';
22     recons = temp(:,1:num_comp)*U(:,1:num_comp)';
23     for j = 1:len(2)
24         for i = 1:len(1)
25             recons(i,j) = recons(i,j) + mean(j);
26         end
27     end
28     err = immse(recons,fatiguedatademo);
29     error(k) = max(err);
30 end
31
32
33 %% To plot the maximum error as function of components considered
34 comp = 1:26;
35 plot(comp,error,'-ro','linewidth',2.0);
36 xlabel('Number of Components');
37 ylabel('Max error in Data')
38 ax = gca;
39 set(ax,'linewidth',1.5);
40 xlim([2 10])
41 axis('square');
42 grid on;
43
44 %% To make scatter plots
45 scatter(c1(1,:),c1(5,:));
```

```

46 xlabel('PC1');
47 ylabel('PC5')
48 ax = gca ;
49 set(ax, 'linewidth',1.5);
50 axis('square');
51 grid on;
52
53
54 %% To make a Pareto plot
55 s = svd(fatiguedatademo','econ') ;
56 pareto(s)
57 xlabel('Components');
58 ylabel('Frequency')
59 ax = gca ;
60 set(ax, 'linewidth',2.0);
61 grid on;
62 legend('Frequency','Cumulative Frequency');

```

The above code was used to create *Pareto plots* and *Scatter Plot* of the principal components calculated for the data.

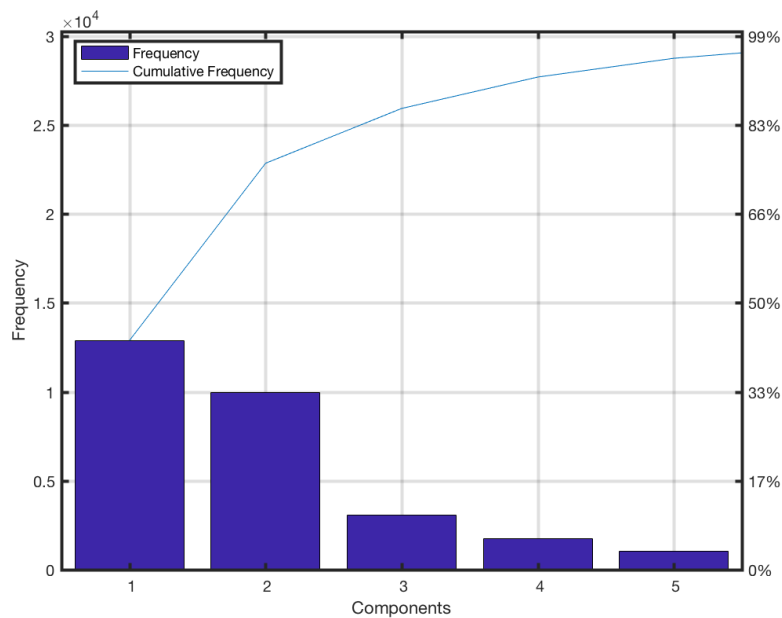
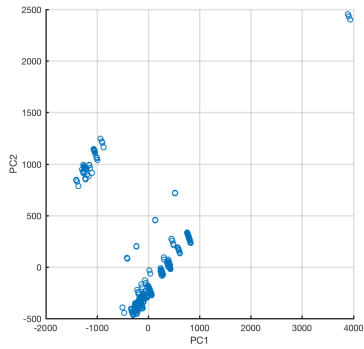
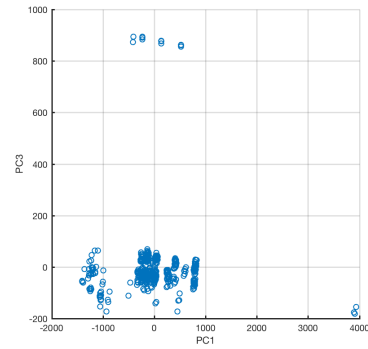


Figure 1: Pareto plot of frequency of Principal Components

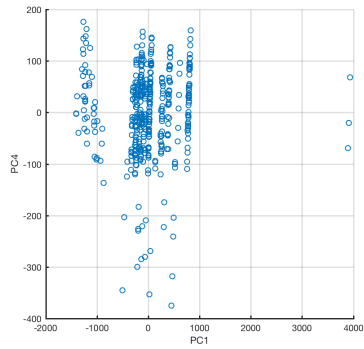
Then we reconstructed the data using the principal components and the maximum error from the original data was calculated as a function of number of components used in Fig. 3



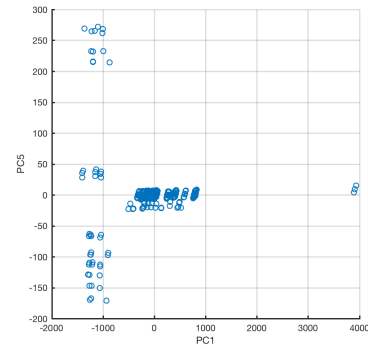
(a) Principal Components 1 and 2



(b) Principal Components 1 and 3



(c) Principal Components 1 and 4



(d) Principal Components 1 and 5

Figure 2: Scatter plots of first 5 principal components

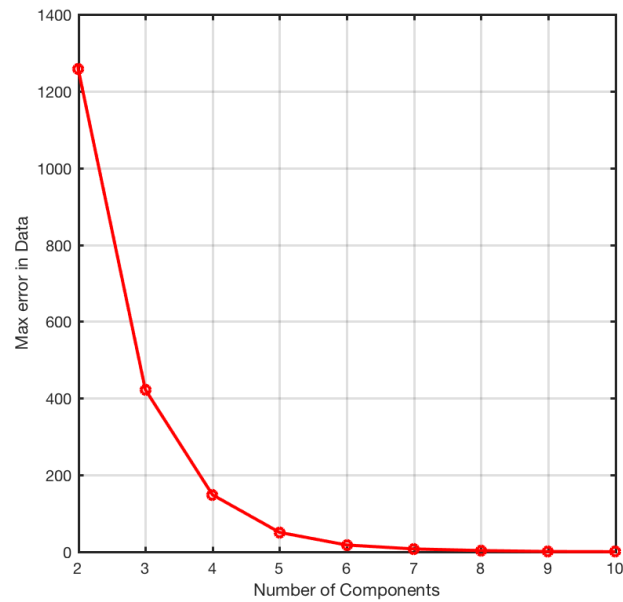


Figure 3: Max. error in reconstructed data

4 Results

We performed **Principal Component Analysis** on the data set provided to us. We made *Pareto plot* and *Scatter plots* on the principal components. We also reconstructed the data using the principal components and compared with the original data to find that 8 components are sufficient to reconstruct the data to reasonable accuracy.