

CSE 601

Data Mining and Bioinformatics

Programming Assignment 1

PCA

Submitted By

Amit There - 50247978 - amitpadm

Sandesh Patnurkar - 50246311 - sandeshv

Principal Component Analysis:

Principal Component Analysis is a method that transfers a set of correlated dimensions into a new set of uncorrelated dimensions. It reduces the number of dimensions of the data. This is especially helpful when the data has a lot of dimensions and you wish to visualize it in a single or a 2 dimensional space or maybe a 3 dimensional space.

Steps:

This is done by following the following sequence of steps in the code:

- 1) Read the file in a matrix.
- 2) We find a mean vector i.e. taking the mean of all the rows.
- 3) Normalize the data using this mean vector.
- 4) Compute the covariance matrix of this normalized vector.
- 5) Finding eigenvalues and eigenvectors of the covariance matrix.
- 6) Sorting on eigenvalues and taking the top d eigenvectors which become the directions with largest variances.
- 7) Rotate the data from original direction to the ones represented by these eigenvectors in 2 dimensional space.
- 8) Plot the rotated data.
- 9) We use the unique diseases from the last column in each file as colours for plotting the data points.

Plots:









