

Principal Components Analysis developed by Karl Pearson in 1901.
PCA combines the attributes and hence reducing

the dimension of a data set.

The basic procedure is as follows.

1. Load the data set. 'A'
2. Normalized the data. Generally zero-mean normalization is done.
$$\text{DataAdjust} = \frac{\sum_{i=1}^n X_i - \bar{x}}{s}$$
3. Find the covariance of the matrix A.
$$\text{COV}(X, Y) = \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) / (n-1)$$
4. Find the Eigen vector (U) and Eigen value (λ) of the covariance matrix.
5. Arrange the eigen value in ~~ascending~~ ^{principal} decreasing order of the data.
6. Arrange the eigen vector according to the eigen value.
7. Keep only those eigen vector for which the eigen value is > 1 .
8. Multiply this reduced eigen vector matrix with the normalized data set, it gives the ~~the~~ principal components of the data set.

Factor Analysis

1. Load the data set
2. Obtain the normalized data set (Y) by subtracting the column wise mean from each sample in the column and divide it by the column wise standard deviation.
3. Calculate the correlation matrix (R) using the relation
$$R = \frac{Y X Y'}{n-1}, \quad n = \text{total no. of samples.}$$
4. Calculate the eigen vector (U) and Eigen value (λ) from R .
5. Rearrange the Eigen values and corresponding Eigen vector in descending order.
6. Calculate the factor loading matrix (A) using $A = U \times \sqrt{\lambda}$
7. Calculate the factor scores (F)

$$B = R^{-1} \times A \quad B = \text{Factor coefficient matrix.}$$
$$\boxed{F = B \times Y.}$$

R^{-1} = inverse of the correlation matrix.

A = factor loading matrix.

Curse of dimensionality

- When dimensionality increases, database becomes increasingly sparse → Information content decreases
- Density and distance between points, which is critical to clustering, outlier analysis becomes less meaningful.
- The possible combinations of subspaces grow exponentially.

Dimensionality Reduction

- Avoid the curse of dimensionality.
- Help eliminate irrelevant features and reduce noise.
- Reduce time and space required in data mining.
- Allow easier visualization.

Dimensionality reduction techniques

DFI

DCF

DWT

PCA

ICA

Lossless - ZIP files

Compression image format, gif, jpeg, png