

PCA

Epoch IIT Hyderabad

Himani Agrawal
MA22BTECH11008

Contents

- Introduction
- Working

Introduction

PCA is a statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components¹. This transformation is defined in such a way that the first principal component has the largest possible variance (that is, accounts for as much of the variability in the data as possible), and each succeeding component in turn has the highest variance possible under the constraint that it is orthogonal to the preceding components.

How does PCA work?

The steps involved in PCA are as follows:

1. **Standardization:** The aim of this step is to standardize the range of the continuous initial variables so that each one of them contributes equally to the analysis.
2. **Compute the Covariance Matrix:** The covariance matrix is a $p \times p$ symmetric matrix (where p is the number of dimensions) that has as entries the covariances associated with all possible pairs of the initial variables. For each pair of initial variables, the

covariance is computed as follows:

$$\text{cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x - \bar{x})(y - \bar{y})$$

$\text{cov}(X, Y) \longrightarrow$ Covariance between X & Y variables

x & $y \longrightarrow$ members of X & Y variables

\bar{x} & $\bar{y} \longrightarrow$ mean of X & Y variables

$n \longrightarrow$ number of members

3. **Compute the Eigenvectors and Eigenvalues of the Covariance Matrix:** Eigenvectors and eigenvalues are the linear algebra concepts that we need to compute from the covariance matrix in order to determine the principal components of the data.
4. **Choose the Principal Components:** To decide which eigenvector(s) can be dropped without losing too much information for the construction of lower-dimensional subspace, we need to inspect the corresponding eigenvalues: The eigenvectors with the lowest eigenvalues bear the least information about the distribution of the data; those are the ones that can be dropped.
5. **Recast the Data Along the Principal Components Axes:** In this final step we will use the eigenvectors to reorient our data from the original axes to the ones represented by the principal components.