

Principal Component Analysis

DR. BHARGAVI R

SCOPE

VIT CHENNAI

Introduction

- Principal Component Analysis (PCA) is a type of unsupervised learning
- A tool used for data visualization or data pre-processing before supervised techniques are applied for high dimensional data.
- Used for deriving a low-dimensional set of features from a large set of interrelated variables.
- Dimensionality reduction : Achieved by transforming to a new set of variables, called principal components (PCs), which are uncorrelated, and which are ordered so that the first few retain most of the variation present in all of the original variables.

Introduction (cont...)

- PCA identifies the patterns in data based on the correlation between the features.
- Finds the directions of maximum variance in high dimensional data and projects the data onto a new subspace with equal or fewer dimensions than the original one.
- The orthogonal axes (PCs) of the new subspace can be interpreted as the directions of maximum variance given the constraint that the new feature axes are orthogonal to each other.

Bhargavi R

PCA - Steps

- Standardize the d-dimensional dataset.
- Compute the *mean* for every dimension of the whole dataset. (i.e compute the mean vector)
- Compute the *covariance matrix* of the whole dataset using the formula:

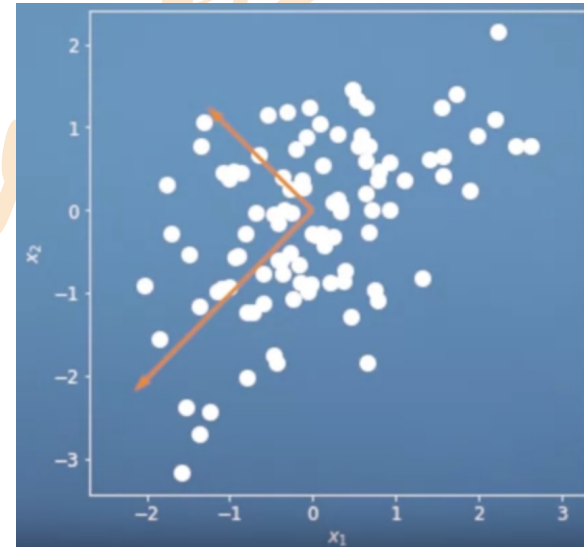
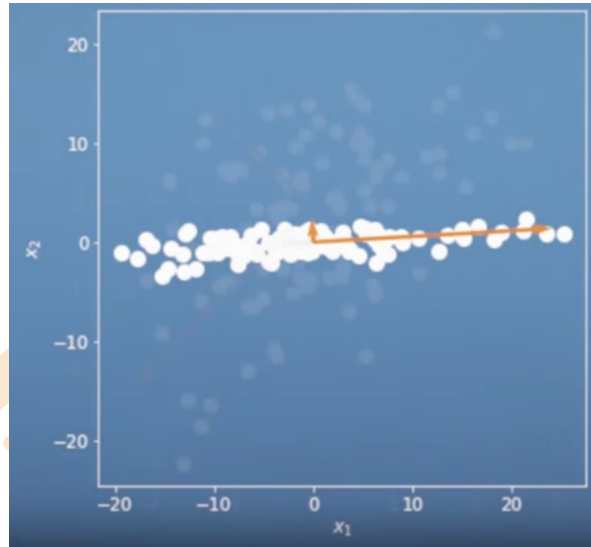
$$\frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \mu)(\mathbf{x}_i - \mu)^T$$

here each \mathbf{x}_i is d dimensional vector & n is the number of instance of data.

- Compute *eigenvectors* and the corresponding *eigenvalues* of the covariance matrix.
- Sort the eigenvectors by decreasing order of eigenvalues and choose k eigenvectors with the largest eigenvalues to form a $d \times k$ dimensional matrix \mathbf{W} .
- Project any data point onto the principal subspace that is spanned by the eigenvectors that belong to the largest eigenvalues (i.e \mathbf{W}).

Data Standardization

Standardizing the data: Data with Zero mean and unit variance



Example

Consider the following dataset

	Sample1	Sample2	Sample3	Sample4
x1	4	8	13	7
x2	11	4	5	15

$$\text{Mean vector} = \begin{matrix} (4 + 8 + 13 + 7)/4 & = & 8 \\ (11 + 4 + 5 + 15)/4 & = & 8.5 \end{matrix}$$

Subtract data from mean of the corresponding feature we get

$$\begin{matrix} -4 & 0 & 5 & -1 \\ 2.5 & -4.5 & -3.5 & 5.5 \end{matrix}$$

Example (cont...)

Find the covariance of the matrix (X-Mean) using the formula $\frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \mu)(\mathbf{x}_i - \mu)^T$

$$\text{Cov}(\mathbf{x}-\text{mean}) = \frac{1}{3} * \begin{bmatrix} -4 & 0 & 5 & -1 \\ 2.5 & -4.5 & -3.5 & 5.5 \end{bmatrix} * \begin{bmatrix} -4 & 0 & 5 & -1 \\ 2.5 & -4.5 & -3.5 & 5.5 \end{bmatrix}^T = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

Now, Calculate the eigen values and eigen vectors of the covariance matrix.

λ is an eigen value for a matrix M if it is a solution of the characteristic equation $\det(\mathbf{M} - \lambda \mathbf{I}) = 0$.

$$\text{i.e } \det \left(\begin{bmatrix} 14 - \lambda & -11 \\ -11 & 23 - \lambda \end{bmatrix} \right) = 0$$

$$\lambda^2 - 37\lambda + 201 = 0$$

Solving the above quadratic equation we get $\lambda_1 = 30.38$, $\lambda_2 = 6.62$

Clearly, the second eigen value is very small compared to the first eigen value.

So, the second eigen vector can be left out if we need only one principal component.

Example (cont...)

Eigen vector corresponding to the greatest eigen value is the principal component for the given data set.

So. we find the eigen vector corresponding to eigen value λ_1 .

We use the following equation to find the eigen vector

$$MV = \lambda V$$

Where M is the Covariance Matrix, V is the Eigen vector & λ is the Eigen value

Substituting the values we get:

$$\begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 30.38 \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$14v_1 - 11v_2 = 30.38v_1$$

$$-11v_1 + 23v_2 = 30.38v_2$$

$$v_1 = -0.67v_2$$

Example (cont...)

One eigen vector could be $\begin{bmatrix} -0.67 \\ 1 \end{bmatrix}$

Now to find the Principal component we need to project the data onto the eigen vector.

i.e find Transpose of Eigen vector x (Feature Vector – Mean Vector)

To find the principal component of the data (4,11)

$$[-0.67 \ 1] * \begin{bmatrix} 4 - 8 \\ 11 - 8.5 \end{bmatrix} = [-0.67 \ 1] * \begin{bmatrix} -4 \\ 2.5 \end{bmatrix} = 5.1$$

Similarly, we we can compute second principal component also.