# Principal Component Analysis (PCA)

# Principal Component Analysis

❖ Principal Component Analysis is an unsupervised learning algorithm that is used for the dimensionality reduction in **machine learning**

❖ It is a statistical process that converts the observations of correlated features into a set of linearly uncorrelated features with the help of orthogonal transformation. These new transformed features are called the Principal Components.

❖ It is one of the popular tools that is used for exploratory data analysis and predictive modeling. It is a technique to draw strong patterns from the given dataset by reducing the variances.
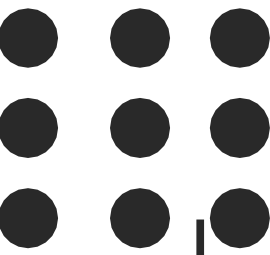
# REAL WORLD APPLICATIONS

- **Image processing**

- Movie recommendation system

- Optimizing the power allocation in various communication channels

# Some common terms used in PCA algorithm:

**Dimensionality:** It is the number of features or variables present in the given dataset. More easily, it is the number of columns present in the dataset.

**Correlation:** It signifies that how strongly two variables are related to each other. Such as if one changes, the other variable also gets changed. The correlation value ranges from -1 to +1. Here, -1 occurs if variables are inversely proportional to each other, and +1 indicates that variables are directly proportional to each other.

# Some common terms used in PCA algorithm:

**Orthogonal:** It defines that variables are not correlated to each other, and hence the correlation between the pair of variables is zero.

**Eigenvectors:** If there is a square matrix M, and a non-zero vector v is given. Then v will be eigenvector if Av is the scalar multiple of v.

**Covariance Matrix:** A matrix containing the covariance between the pair of variables is called the Covariance Matrix.

# Steps for PCA Algorithm

**1. Getting the dataset :**

Firstly, we need to take the input dataset and divide it into two subparts X and Y, where X is the training set, and Y is the validation set.

**2. Representing data into a structure :**

Now we will represent our dataset into a structure. Such as we will represent the two-dimensional matrix of independent variable X. Here each row corresponds to the data items, and the column corresponds to the Features. The number of columns is the dimensions of the dataset.

## 3. Standardizing the data

In this step, we will standardize our dataset. Such as in a particular column, the features with high variance are more important compared to the features with lower variance.
If the importance of features is independent of the variance of the feature, then we will divide each data item in a column with the standard deviation of the column. Here we will name the matrix as Z.

## 4. Calculating the Covariance of Z

To calculate the covariance of Z, we will take the matrix Z, and will transpose it. After transpose, we will multiply it by Z. The output matrix will be the Covariance matrix of Z.

## 5. Calculating the Eigen Values and Eigen Vectors

Now we need to calculate the eigenvalues and eigenvectors for the resultant covariance matrix Z.
Eigenvectors or the covariance matrix are the directions of the axes with high information. And the coefficients of these eigenvectors are defined as the eigenvalues.

## 6. Sorting the Eigen Vectors

In this step, we will take all the eigenvalues and will sort them in decreasing order, which means from largest to smallest.
And simultaneously sort the eigenvectors accordingly in matrix P of eigenvalues. The resultant matrix will be named as P*.

**7. Calculating the new features Or Principal Components**

Here we will calculate the new features. To do this, we will multiply the P* matrix to the Z. In the resultant matrix Z*, each observation is the linear combination of original features. Each column of the Z* matrix is independent of each other.

**8. Remove less or unimportant features from the new dataset**

The new feature set has occurred, so we will decide here what to keep and what to remove.
It means, we will only keep the relevant or important features in the new dataset, and unimportant features will be removed out.

# Properties of Principal components

❖ The principal components are must be the linear combinations of the original variables, the weights vector in this combination is the eigenvector found which will satisfies the principle of least squares.

❖ Principal components are orthogonal in nature which means linearly dependent.

❖ The variation or spread of the data in the principal components decreases as we move from the first principal component to the last.

# Advantages of Principal Component Analysis

1.It removes Features with correlation.

2.It improves performance of the algorithm.

3.It also reduces overfitting.

4.Visualization of data is improved in this method.

# Disadvantages of Principal component Analysis

1.In this method, independent variables become less interpretable.

2.Before beginning PCA, we must perform Data standardization.

3.There will be information loss by using this method.