

⇒ PCA & Principal Component Analysis-

it is used for to reduced the dimensionality.

it is used in the unsupervised learning Algorithm.

⇒ LDA (Linear discriminant Analysis) Algorithm.

which is used in the supervised Learning Algorithm to reduce the dimension.

→ Definition- Principal component Analysis (PCA) is a dimensionality reduction technique used in unsupervised Learning to transform a dataset with many variables (features) into a smaller set that still contains most of the essential information.

e.g. we have a dataset.

and if we have to calculate the price of the house. then what we see mostly ?? like the size of the plot, or number of rooms, or number of washrooms or number of floors like that. we take the multiple features.

so if we have less features we can get the good accuracy but if we have many features then there complexity increase there, and at that time accuracy become decrease. and prediction will no come good.

so that's the reason we have to reduce the dimension.

→ steps of principal component Analysis (PCA)

1 → Calculate Mean of every feature.

2 Calculate the covariance Matrix.

$$\text{Cov} = \begin{bmatrix} \text{Var}(n_1) & \text{cov}(n_1, n_2) \\ \text{cov}(n_1, n_2) & \text{Var}(n_2) \end{bmatrix}$$

example we have data set -

Size (in cm)

n_1

5

3

10

2

20 → sum

4 → Total Num

Mean of n_1 = 5

No of rooms

n_2

10

7

15

4

36 → sum of Number

4 → Total no of n_2

Mean of n_2 = 9

so

$$\bar{n}_1 = 5$$

$$\bar{n}_2 = 9$$

step 2-

Calculate Covariance Matrix

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} 12.67 & 16.33 \\ 16.33 & 22 \end{bmatrix}$$

so how can we calculate the value of σ_{11} ?

we have to calculate variance.

$$\sigma_{11} = \frac{(5-5)^2 + (3-5)^2 + (10-5)^2 + (2-5)^2}{N-1} = \frac{50}{3}$$

σ_{11} → mean of n_1
value of n_1 → mean of n_1
first value → value of n_1
second value

$$N-1$$

(where N = total number of values in n_1 , which is 4, then $n-1 = 4-1 = 3$)

we will get $[12.67] \rightarrow \sigma_{11}$

same way find the var(n_2) → σ_{22} .

you will get 22

Note \Rightarrow What is variance? Variance is a number that tells you how spread your data is.

\rightarrow Imagine you have test scores from two different classrooms - class A and B. Class A's scores 71, 72, 73, 74, 75

Class B's scores 50, 60, 73, 86, 96

\rightarrow The average (mean) Test score for both class is 73.

But the two sets of score are very different!

- In class A, every student scored very close to the average (mean) 73, The scores are bunched together.
- In class B, some students did very poorly and some very well, The scores are all over the place, They are spread out.
- Variance is the number that captures this "spread".
- Class A would have a low variance because the scores are close together.
- Class B would have a high variance because the scores are far apart.

→ Now let's see how can we calculate σ_{12} ? $\text{cov}(n_2, n_1)$ and $\text{cov}(n_1, n_2)$

→ So for that,

$$\frac{(n_1 - \bar{n}_1)(n_2 - \bar{n}_2)}{n-1}$$

$$\frac{(5-5)(10-9) + (3-5)(7-9) + (10-5)(2-5)(4-9)}{3}$$

we get

$$\sigma_{12} = 16.33 \rightarrow \text{both are same}$$

$$\sigma_{21} = 16.33 \rightarrow \text{because it is in the multiplication.}$$

\Rightarrow

Step 3.8 Calculate the Eigenvalues

$$|S - \lambda I| = 0$$

\downarrow
Covariance
Matrix.

$$\begin{vmatrix} 12.67 - \lambda & 16.33 \\ 16.33 & 22 - \lambda \end{vmatrix} = 0$$

and after finding the determinants of it we will get.

$$\lambda_1 = 32.3$$

$$\lambda_2 = -0.4$$

we have two Eigen values because we worked on two features n_1 and n_2 .

\Rightarrow so from the both values we have to consider the largest value for calculate the next step which is Eigen vectors which is largest value is $\lambda_1 = 32.3$.

Step 4:- Calculate the Eigen vectors

$$\begin{bmatrix} 12.67 - 1 & 16.33 \\ 16.33 & 22 - 1 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} -1.3256 \\ 1 \end{bmatrix}$$

here v_1 and v_2 are the vectors (original data)

Step 5: Project Data onto principal components

$$\begin{bmatrix} 5 & 10 \\ n_1 & n_2 \end{bmatrix} \begin{bmatrix} -1.3256 \\ 1 \end{bmatrix} \begin{matrix} u_1 \\ u_2 \end{matrix}$$

1x2 columns
2x1

we will get 1x1 Matrix

$$[5 \times -1.3256 + 10 \times 1] = 3.03 \rightarrow P.C$$

same for others

3

n_1 n_2 P.C

5	10	3.03
3	7	3
10	15	10.8
2	4	1.03

so it means we have reduce the 2 dimensions into one, n_1 and n_2 to P.C.

Now we can use just P.C feature to predict the data.