# Dimension Reduction Method

**\* Dimentionality Reduction:**

→ It is a process of reducing no. of input variables or features in a dataset while retaining meaning information.

→ It is important step in data mining because in real-time application, dataset contains many attributes not all of which are useful or relevant

**\* Multi-colinearity:-**

It occures when two or more independent variables in a dataset are highly correlated with each other.

Ex - Predicting house price using Area, types of, locality. Among these area and types are highly correlated and processes multicolinearity.

Problems:

(i) Unstable coefficient in regression Model.

(ii) Reduces Model interpreetability.

(iii) Decrease Model performance.

Handle:

(i) Use PCA (principle component Analysis).

→ Other two methods are:-

2mark (a) Ridge (b) Lano

**\* Need of Dimentionality Reduction:-**

(i) To remove irrelevant and redundant features.

(ii) To improve model performance.

(iii) To reduce storage and memory requirement.

(iv) To improve visualization and interpretation.

**\* Principle Component Analysis:-**

Steps:-

(1) Standarcised the data using Z-score Normalization.

$$Z\text{-}Score = \frac{x_i - \overline{x_i}}{S_i}$$

where $\overline{x_i}$, mean of $x_i$

$S_i$, standared deviation of $x_i$

(2) Compute co-varriance matrix

(3) Compute eigen value & eigen vector.

(4) Compute PCA.

(Q) Suppose we have variables exam1, exam2 of steedents.

| Student | Exam1 | Exam2 |
|---------|-------|-------|
| A | 90 | 85 |
| B | 70 | 65 |
| C | 80 | 78 |
| D | 65 | 60 |
| E | 95 | 92 |

Apply PCA to compute principle component.

Sol$^n$ **Step-1**

Mean of Exam1 $= \dfrac{90+70+80+65+95}{5} = \cancel{300}\ 80$

Mean of Exam2 $= \dfrac{85+65+78+60+92}{5} = 76$

S.D of Exam1 $= \sqrt{\dfrac{\sum\limits_{i=1}^{n}(x_i - \overline{x_i})^2}{n-1}}$

$= \sqrt{\dfrac{(90-80)^2 + (70-80)^2 + (80-80)^2 + (65-80)^2 + (95-80)^2}{4}}$

$= \sqrt{\dfrac{(10)^2 + (-10)^2 + (0)^2 + (-15)^2 + (15)^2}{4}}$

$= \sqrt{\dfrac{100 + 100 + 225 + 225}{4}} = \cancel{100.6}\ 12.748$

$= \cancel{12.748}$

$$S.D \text{ of } Exam2 = \sqrt{\frac{(85-76)^2 + (65-76)^2 + (78-76)^2 + (60-76)^2 + (92-76)^2}{4}}$$

$$= 13.398$$

Z-score$_A$ (Exam 1) $= \dfrac{90-80}{12.748} = \dfrac{10}{12.748} = 0.78$

Z-score$_B$ (Exam 1) $= \dfrac{70-80}{12.748} = \dfrac{-10}{12.748} = -0.78$

Z-score$_C$ (Exam 1) $= \dfrac{80-80}{12.748} = \dfrac{0}{12.748} = 0$

Z-score$_D$ (Exam 1) $= \dfrac{65-80}{12.748} = \dfrac{-15}{12.748} = -1.17$

Z-score$_E$ (Exam 1) $= \dfrac{95-80}{12.748} = \dfrac{15}{12.748} = 1.17$

Z-score$_A$ (Exam 2) $= \dfrac{85-76}{13.398} = 0.671$

Z-score$_B$ (Exam 2) $= \dfrac{65-76}{13.398} = -0.821$

Z-score$_C$ (Exam 2) $= \dfrac{78-76}{13.398} = 0.149$

Z-score$_D$ (Exam 2) $= \dfrac{60-76}{13.398} = -1.194$

Z-score$_E$ (Exam 2) $= \dfrac{92-76}{13.398} = 1.194$

Step-2 Compute covarreiance matrix for standardize root data.
Covarreiance matrix!-
   It is a square matrix that contain covarrciance
between each paire in the dataset.

$$Cov(Exam1, Exam2) = \frac{\sum_{i=1}^{n} x_i^2 (exam1)}{n-1}$$

$$= \frac{(0.784)^2 + (0.784)^2 + (0)^2 + (-1.177)^2 + (1.177)^2}{4}$$

$$= \frac{0.614 + 0.614 + 0 + 1.385 + 1.385}{4} = 0.995 \approx 1.$$

$$\text{Cov (Exam2)} = \frac{\sum_{i=1}^{n} x_i^2 \, (\text{exam2})}{n-1}$$

$$= \frac{(0.672)^2 + (-0.821)^2 + (0.149)^2 + (-1.194)^2 + (1.194)^2}{4}$$

$$\approx 1$$

$$\text{Cov (Exam1, Exam2)} = \frac{\sum_{i=1}^{n} x_i(\text{exam1}) \cdot x_i(\text{exam2})}{n-1}$$

$$= \frac{\begin{matrix}(0.784)(0.671) + (-0.784)(-0.821) + (0)(0.149) + \\ (-1.17)(-1.194) + (1.17)(1.194)\end{matrix}}{4}$$

$$= \frac{0.526 + 0.643 + 0 + 1.396 + 1.396}{4} = 0.99$$

$$S = \begin{bmatrix} 1 & 0.99 \\ 0.99 & 1 \end{bmatrix} \quad \longrightarrow \quad \text{①}$$

Step 3

Eigen value $= |S - \lambda I| = 0$.

$$\lambda I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \lambda = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$S - \lambda I = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0.99 \\ 0.99 & 1 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = \begin{bmatrix} 1-\lambda & 0.99 \\ 0.99 & 1-\lambda \end{bmatrix} \quad \longrightarrow \quad \text{②}$$

$$\Rightarrow \lambda_1 = 1.995, \quad \lambda_2 = 0.005.$$

for eigen vector, $(S - \lambda I)\vec{V} = 0$

$$\Rightarrow \lambda_1 = 1.995$$

Now,

$$\vec{V} = \begin{bmatrix} x \\ y \end{bmatrix} \Rightarrow \begin{bmatrix} 1 - 1.995 & 0.995 \\ 0.995 & 1 - 1.995 \end{bmatrix} \times \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} -0.995 & 0.995 \\ 0.995 & -0.995 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} -0.995x + 0.995y \\ 0.995x - 0.995y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Now,

put $y = x$,

$$\vec{V} = \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x \\ x \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Then,

$$\|\vec{V}\| = \sqrt{1^2 + 1^2} = \sqrt{2} = 1.414$$

$$\vec{V_1} = \frac{1}{1.414} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.707 \\ 0.707 \end{bmatrix}$$

$$\vec{V_2} = \begin{bmatrix} 0.707 \\ -0.707 \end{bmatrix}$$

step 4

Select principle component.

Variance captured by PC1 $= \dfrac{1.995}{\lambda_1 + \lambda_2} = \dfrac{1.995}{1.995 + 0.005}$

$$= 0.9975 \times 100$$

$$= 99.75\%$$

$$PC2 = \frac{0.005}{2} = 0.25\%$$

∴ PC1 captures all the structures in the dataset than PC2.

**Step-5** Transform the data using PC:

~~PCA~~ $\vec{V_1} = \begin{bmatrix} 0.707 \\ 0.707 \end{bmatrix}$ PC

$PC1_A = 0.707 \times$ z-score (Exam1) $+ 0.707 \times$ z-score (Exam2).

$\quad = 0.707 \times 0.78 + 0.707 \times 0.671$

$\quad = 0.551 + 0.474 = 1.025$

$PC1_B =$