

Data Science

Q) PCA for

x	2.5	0.5	2.2	1.9	3.1	2.3	2.0	1.0	1.5	1.1
y	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9

n = 10

x	y	$x_i - \bar{x}$ (A)	$y_i - \bar{y}$ (B)	AB	A ²	B ²
2.5	2.4	0.69	0.49	0.3381	0.4761	0.2401
0.5	0.7	-1.31	-1.21	1.5851	1.7161	1.4641
2.2	2.9	0.39	0.99	0.3861	0.1521	0.9801
1.9	2.2	0.09	0.29	0.0261	0.0081	0.0841
3.1	3.0	1.29	1.09	1.4061	1.6641	1.1881
2.3	2.7	0.49	0.79	0.3871	0.2401	0.6241
2.0	1.6	0.19	-0.31	-0.4589	0.0361	0.0961
1.0	1.1	-0.81	-0.81	0.6561	0.6561	0.6561
1.5	1.6	-0.31	-0.31	0.0961	0.0961	0.0961
1.1	0.9	-0.71	-1.01	0.7171	0.5041	1.0201
$\bar{x} = 1.81 \quad \bar{y} = 1.91$				5.539	5.549	5.449

Covariance matrix =
$$\begin{bmatrix} \text{Cov}(x, x) & \text{Cov}(x, y) \\ \text{Cov}(y, x) & \text{Cov}(y, y) \end{bmatrix}$$
 ... here $\text{Cov}(x, y) = \text{Cov}(y, x)$

where; $\text{Cov}(x, x) = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$; $\text{Cov}(x, y) = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1}$

⇒ Covariance matrix =
$$\frac{1}{9} \begin{bmatrix} 5.549 & 5.539 \\ 5.539 & 6.449 \end{bmatrix} = \begin{bmatrix} 0.6166 & 0.6154 \\ 0.6154 & 0.7166 \end{bmatrix}$$

now, we need to find eigen values & eigen vectors for the above covariance matrix.

$$\Rightarrow |A - \lambda I| = 0$$

$$\Rightarrow \begin{bmatrix} 0.6166 & 0.6154 \\ 0.6154 & 0.7166 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = 0$$

$$\Rightarrow \begin{bmatrix} (0.6166 - \lambda) & 0.6154 \\ 0.6154 & (0.7166 - \lambda) \end{bmatrix} = 0 \quad \text{--- eq (1)}$$

By finding det for above matrix, we will get an equation, which can be solved to get the eigen values:

$$\Rightarrow (0.6166 - \lambda)(0.7166 - \lambda) - (0.6154)^2 = 0$$

$$\Rightarrow 0.441 - 0.7166\lambda - 0.6166\lambda + \lambda^2 - 0.378 = 0$$

$$\Rightarrow \lambda^2 - 1.3332\lambda + 0.063 = 0$$

now, we will find roots for above equation: $\lambda = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$

$$\Rightarrow \lambda_1 = 1.283 \quad \& \quad \lambda_2 = 0.049$$

now, putting λ_1 in eq. (1) for eigen vectors

$$\begin{bmatrix} (0.6166 - 1.283) & 0.6154 \\ 0.6154 & (0.7166 - 1.283) \end{bmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

(or)

eigen vectors can be found by using: $CV = \lambda V$

$$\Rightarrow \begin{bmatrix} 0.6166 & 0.6154 \\ 0.6154 & 0.7166 \end{bmatrix} \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = 1.283 \begin{bmatrix} x_1 \\ y_1 \end{bmatrix}$$

$$\Rightarrow 0.6166x_1 + 0.6154y_1 = 1.283x_1 \Rightarrow x_1 = 0.923y_1$$

$$\Rightarrow 0.6154x_1 + 0.7166y_1 = 1.283y_1$$

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} 0.923 \\ 1 \end{bmatrix}$$

$$\therefore \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} \frac{x_1}{\sqrt{x_1^2 + y_1^2}} \\ \frac{y_1}{\sqrt{x_1^2 + y_1^2}} \end{bmatrix} = \begin{bmatrix} \frac{0.923}{\sqrt{1.851}} \\ \frac{1}{\sqrt{1.851}} \end{bmatrix} = \begin{bmatrix} \frac{0.923}{1.36} \\ \frac{1}{1.36} \end{bmatrix} = \begin{bmatrix} 0.678 \\ 0.735 \end{bmatrix}$$

for $\lambda = 0.049$ $\therefore CV = \lambda V$

$$\begin{bmatrix} 0.6166 & 0.6154 \\ 0.6154 & 0.7166 \end{bmatrix} \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = 0.049 \begin{bmatrix} x_1 \\ y_1 \end{bmatrix}$$

$$\Rightarrow \begin{aligned} 0.6166x_1 + 0.6154y_1 &= 0.049x_1 \\ 0.6154x_1 + 0.7166y_1 &= 0.049y_1 \end{aligned} \Rightarrow x_1 = \frac{-0.6154}{0.5676} y_1$$

$$\Rightarrow x_1 = -1.084 y_1$$

$$\therefore \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} \frac{-1.084}{\sqrt{2.175}} \\ \frac{1}{\sqrt{2.175}} \end{bmatrix} = \begin{bmatrix} \frac{-1.084}{1.474} \\ \frac{1}{1.474} \end{bmatrix} = \begin{bmatrix} -0.735 \\ 0.678 \end{bmatrix}$$

SVD: $A_{m \times n} \rightarrow \text{rank}(A)$

$$\Sigma_{m \times n}; U_{m \times m}; V_{n \times n}$$

$$\Rightarrow A = U \Sigma V^T$$

- from given "A" find $A^T A$ & find eigen value for that matrix
 - and also find σ values [like $\sigma_1 = \sqrt{\lambda_1}$; $\sigma_2 = \sqrt{\lambda_2}$, etc...]
 - from eigen values find eigen vectors and place beside each other by normalizing them to get "V"
 - find Σ , as it is a diagonal matrix of values " σ " & make its rank equal to A i.e. $\text{Rank}(\Sigma) = \text{Rank}(A)$.
 - now $U_1 = \frac{1}{\sigma_1} A V_1$ & $U_2 = \frac{1}{\sigma_2} A V_2 \Rightarrow U = [U_1, U_2]$
 - $U = \text{normalized eigen vector matrix } A A^T$
 - $V^T = \text{normalized eigen vector matrix } A^T A$
- $$A = \underbrace{U}_{\substack{\sqrt{A^T A} \\ A A^T}} \underbrace{\Sigma}_{\substack{\downarrow \\ \text{diagonal} \\ \text{matrix}}} \underbrace{V^T}_{\substack{\uparrow \\ A^T A}}$$

Statistics is the study of the collection, analysis, interpretation, presentation and organization of data.

eg:- mean of marks obtained by 50 students in a class.

Basics of statistics include the measure of central tendency (mean, median, mode) and the measure of dispersion (variance and standard deviation).

Types of Statistics:-

i) Descriptive Statistics :-

These statistics are used to summarize and describe a data set. Histograms, pie charts, bar and scatter plots are common ways to summarize data and present it in table and graphs.

Here, the summarization is done using:-

- Measures of central tendency : mean, median, mode
- Measures of dispersion : range, variance, standard deviation.
- Measures of shape : skewness, kurtosis.

For instance, if we have a data set containing the weights of 20 people, we could use descriptive statistics to summarize and describe the data, such as calculating the mean weight, the range of weights, or the skewness of the weight distribution.

ii) Inferential Statistics:-

These statistics are used to draw conclusions or make predictions about a population based on a sample of data.

Here, the prediction is done using (t-tests, ANOVA, chi-square test & linear regression)

- Hypothesis testing: testing whether a sample mean is significantly different from population mean
- Confidence intervals: estimating the range of values within which a population parameter (like mean) likely to fall.
- Regression analysis: analyzing the relationship between variables and predicting values based on that relationship.

For example, we might use inferential statistics to test whether a new drug reduces cholesterol levels in a population. We would take a sample of people and randomly assign them to a treatment group or a control group, and then use statistical tests to determine whether there is a significant difference in cholesterol levels between the two groups.

Probability :-

Probability is a mathematical concept that predicts how likely events are to occur. The probability values are expressed between 0 & 1. The definition of probability is the degree to which something is likely to occur.

There are mainly two types of probability distributions:

i) Discrete probability :-

These distributions are associated with random variables that can only take on a finite or countable number of values, such as the outcome of a coin toss or the number of defective items in a batch of products.

ii) Continuous probability :-

These distributions on other hand are associated with random variables that can take on any value within a range, such as the height or weight of individuals in a population.

Examples of common probability distribution include the normal distribution (also known as Gaussian distribution), the binomial distribution, the poisson distribution and the exponential distribution.

Probability distributions are widely used in fields such as statistics, physics, engineering, finance, and many other areas of science and industry to model and analyze a wide range of phenomena.

SVD (singular value decomposition) :-

$$A_{m \times n} = U_{m \times m} \cdot \Sigma_{m \times n} \cdot V^T_{n \times n}$$

rotation stretching rotation

here U & V are orthogonal i.e. $U \cdot U^T = I$
 $V \cdot V^T = I$
 & Σ is a rectangular matrix with singular values

$$A^T \cdot A = (V \cdot \Sigma^T \cdot U^T) (U \cdot \Sigma \cdot V^T)$$

$$= V \cdot \Sigma^T \cdot \Sigma \cdot V^T$$

rotation: $\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$ → unitary transformation

$$A \cdot A^T = (U \cdot \Sigma \cdot V^T) (V \cdot \Sigma^T \cdot U^T)$$

$$= U \cdot \Sigma \cdot \Sigma^T \cdot U^T$$

stretching: $\begin{pmatrix} \alpha & 0 \\ 0 & \alpha \end{pmatrix}$

eg: $A = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix}_{2 \times 3}$

$A^T \cdot A = \Sigma$
 $A \cdot A^T = U$

$$A^T = \begin{bmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{bmatrix}_{3 \times 2}$$

$$A \cdot A^T = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix} \begin{bmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 11 & 5 \\ 5 & 8 \end{bmatrix}$$

$\Rightarrow (A - \lambda I) x = 0$

$\Rightarrow \begin{vmatrix} 11-\lambda & 5 \\ 5 & 8-\lambda \end{vmatrix} = 0$

$(11-\lambda)^2 - 5^2 = 0 \Rightarrow (11-\lambda)^2 = 25$
 $\Rightarrow 11-\lambda = \pm 5$
 $\Rightarrow \lambda = 10; \lambda = 12$

Put λ values in above equation.

$\begin{vmatrix} 11-10 & 5 \\ 5 & 8-10 \end{vmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0 \Rightarrow \begin{vmatrix} 1 & 5 \\ 5 & -2 \end{vmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$

$\begin{vmatrix} 11-12 & 5 \\ 5 & 8-12 \end{vmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0 \Rightarrow \begin{vmatrix} -1 & 5 \\ 5 & -4 \end{vmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$

we need to write in descending order so,

now calculating orthogonal for above

$x_1 + x_2 = 0$
 $x_1 = -x_2$
 $x_1 + x_2 = 0$
 $x_1 - x_2 = 0$
 $\therefore x_1 = x_2$

$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \xrightarrow{\sqrt{2}} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}$

diagonalized $\rightarrow \begin{pmatrix} 0 & 0 & 5/6 \\ 0 & 5/6 & 0 \end{pmatrix}$

now for $V: A^T A$.

$$\Rightarrow \begin{vmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{vmatrix} \Rightarrow \begin{vmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{vmatrix} \Rightarrow \begin{vmatrix} 10 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{vmatrix}$$

$$\lambda^3 - S_1 \lambda^2 + S_2 \lambda - S_3 = 0$$

$$S_1 = 22$$

$$S_2 = 100 + (4) + 16 = 120$$

$$S_3 = 10(4) - 0 + 2(-20) = 0$$

$$\Rightarrow \lambda^3 - 22\lambda^2 + 120\lambda = 0$$

$$\Rightarrow \lambda(\lambda^2 - 22\lambda + 120) = 0$$

$$\Rightarrow \lambda^2 - 22\lambda + 120 = 0$$

$$\therefore \lambda = 0; \lambda = 10; \lambda = 12 \Rightarrow \lambda_1 = 12; \lambda_2 = 10; \lambda_3 = 0$$

$$\Rightarrow \lambda(\lambda - 12) - 10(\lambda - 12) = 0 \Rightarrow (\lambda - 10)(\lambda - 12)$$

by cross rule:-

$$\therefore \text{for } \lambda = 12: \begin{vmatrix} (0-12) & 0 & 2 \\ 0 & (10-12) & 4 \\ 2 & 4 & (2-12) \end{vmatrix} \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix}$$

$$\Rightarrow \frac{x_1}{4} = \frac{-x_2}{-8} = \frac{x_3}{4}$$

$$\begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$$

$$\text{for } \lambda = 10: \begin{vmatrix} 0 & 0 & 2 \\ 0 & 0 & 4 \\ 2 & 4 & 2 \end{vmatrix} \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} \Rightarrow \frac{x_1}{-16} = \frac{-x_2}{-8} = \frac{x_3}{0} \Rightarrow \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} = \begin{bmatrix} 2 \\ -1 \\ 0 \end{bmatrix}$$

$$\text{for } \lambda = 0: \begin{vmatrix} 0 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{vmatrix} \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} \Rightarrow \frac{x_1}{4} = \frac{-x_2}{-8} = \frac{x_3}{-20} \Rightarrow \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} = \begin{bmatrix} 1 \\ 2 \\ -5 \end{bmatrix}$$

$$V = \begin{bmatrix} 1 & 2 & 1 \\ 2 & -1 & 2 \\ 1 & 0 & -5 \end{bmatrix}$$

$$\Rightarrow \text{orthogonal } V^T = \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{6}} & \frac{1}{\sqrt{6}} \\ \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{5}} & 0 \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{30}} & \frac{5}{\sqrt{30}} \end{bmatrix}$$

$$\text{orthogonal } V = \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{6}} & \frac{1}{\sqrt{6}} \\ \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{5}} & 0 \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{30}} & \frac{5}{\sqrt{30}} \end{bmatrix}$$

Σ = diagonal matrix with same dimensionality as A

$$\Rightarrow \begin{bmatrix} \sqrt{\lambda_1} & 0 & 0 \\ 0 & \sqrt{\lambda_2} & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

here λ_1 & λ_2 are 12 & 10 as

they are common in both $A A^T$; $A^T A$ and should be in descending order

$$\therefore A = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{6}} & \frac{1}{\sqrt{6}} \\ \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{5}} & 0 \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{30}} & \frac{5}{\sqrt{30}} \end{bmatrix}$$

Data preprocessing :-

It is the process of transforming raw data into a format that can be easily analyzed by ML algorithms. It involves several steps, including data cleaning, integration, transformation and reduction.

i) Data cleaning :-

- process of identifying and removing (or) correcting errors, inconsistencies and the missing values from data.
- Techniques used for data cleaning includes imputation, deduplication, outlier removal and handling missing values.
- Example: filling in missing values with estimated values when there is missing data in the 'age' column of dataset, this is called imputation.

ii) Data integration :-

- Process of combining data from multiple sources and resolving inconsistencies or conflicts.
- Techniques used includes data fusion, data alignment and data merging.
- Example: merging of customer data based on common identifier, such as customer ID and then align the information.

iii) Data transformation :-

- Process of converting data into a format suitable for machine learning algorithms.
- Techniques used includes normalization, feature scaling, encoding & data discretization.
- Example: Converting categorical variable like colour into numerical value by encoding enables the ML algorithm to use in it directly.

iv) Data Reduction :-

- Process of reducing the amount of data while retaining as much relevant info as possible.
- Techniques used includes feature selection, dimensionality reduction and data sampling.
- Example: Selecting a subset of features that are most relevant to the analysis by dimensionality reduction to increase computation speed of the model.

Student T-test:-

Student T-test is a statistical hypothesis test that is used to compare the mean of two independent group of data. The t-test calculates a t-value, which is a ratio of the difference between the two means to the standard error of the difference.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{S_{\bar{x}}}$$

$$\text{where } S_{\bar{x}} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}}$$

where \bar{x}_1 = mean of first set of values,

\bar{x}_2 = mean of second set of values,

s_1 = standard deviation of first set of values,

" second " " "

n_1 = Total no of values in first set,

n_2 = " " second "