



## Unit -3

### Introduction to Statistical Learning Theory

Features are the variables that can be defined and observed. Feature extraction is the process of identifying and selecting the most important characteristics from a dataset without losing vital information.

Feature extraction essentially is the process of converting raw data into numerical features that can be processed while preserving the information in the original data-set.

- Feature extraction is helpful when training a machine learning model. It leads to:
- An improvement in model accuracy
  - A boost in training speed
  - A reduction in cost of overfitting
  - Better data interpretation

## Feature extraction techniques:

- Principle Components Analysis (PCA)
- Independent Component Analysis (ICA)
- Linear Discriminant Analysis (LDA)
- Locally Linear Embedding (LLE)
- etc.

Challenges in Feature Extraction:

- ↳ Choosing the right method
- ↳ Loss of information
- ↳ Computational complexity

Dimensionality reduction :-

In machine learning model, the higher the number of features, the harder it gets to visualize the training set and then work on it. Sometimes, many of these features are correlated or redundant. This is where dimensionality reduction ~~as~~ algorithms come into play.

Dimensionality reduction can be done in two ways:

- Feature extraction
- Feature selection

The feature extraction technique gives us new features which are a linear combination of the existing features. The new set of features will have different values as compared to the original feature value. The main aim is that fewer features will be required to capture the same info.



## Principal Component Analysis

PCA used in data analysis, particularly for reducing the dimensionality of datasets while preserving crucial information. It does this by transforming the original variables into a set of new, uncorrelated variables called principal components.

### PCA key aspects :-

- ↳ Dimensionality Reduction
- ↳ Linear transformation
- ↳ feature selection
- ↳ Data compression
- ↳ Matrix requirement PCA work with symmetric correlation or covariance matrices and requires numeric, standardized data.
- ↳ Eigen values and Eigen vectors Eigen values represent variance magnitude, and Eigen vectors indicate variance direction.
- ⇒ Limitation: It may loss information,

## How PCA work?

1. Standardize the data

If the features of dataset are on different scales, it's essential to standardize them (subtract the mean and divide by the Standard deviation).

2. Compute the covariance matrix.

Calculate the covariance matrix for the standardized dataset.

3. Compute Eigen vectors and Eigenvalues of covariance matrix.

4. Sort Eigenvectors by Eigenvalues in descending order

5. Choose principal components  
Select the top  $K$  eigenvectors (principal components) where  $K$  is the desired dimensionality of the reduced dataset.

6. Transform the Data

Multiply the original standardized data by the selected Principal Components to obtain the new, lower-dimensional representation of the data.



# POORNIMA

COLLEGE OF ENGINEERING

## DETAILED LECTURE NOTES

PAGE NO. ....

### Singular Value Decomposition :-

SVD is primarily used for dimensionality reduction, information extractors, and noise reduction.

Importance in Data Science :

- Dimensionality Reduction
- Noise Reduction
- Recommendation System

SVD is a method for decomposing a matrix.

For a given matrix A, it can be decomposed into three matrices U, Σ and V<sup>T</sup> such that

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

Given a matrix A of dimension m × n, the SVD decomposes it into three matrices:

1. U : An m × m orthogonal matrix, called the left singular matrix.
2. Σ : An m × n diagonal matrix, with non-negative real numbers on the diagonal. The values on the diagonal are the "singular values" and are usually ordered in decreasing order.

3. VT: An  $n \times n$  orthogonal matrix, where  
 $V^T$  is the transpose of  $V$ , called the right singular matrix.

So, the decomposition can be represented as:

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

The columns of  $U$  are called the left singular vectors, the columns of  $\Sigma$  are called the right singular vectors, and the values in  $\Sigma$  are the singular values of  $A$ .

Eg:-

Compute the SVD for small matrix

$$A = \begin{bmatrix} 4 & 0 \\ 0 & 3 \end{bmatrix}$$

(using Numpy we get this way)

$$\rightarrow U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 4 & 0 \\ 0 & 3 \end{bmatrix}$$

$$V^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$A$  is already a diagonal matrix, so SVD just gives us back the matrix  $A$  itself in  $\Sigma$  and  $U$ ,  $V^T$  are identity matrix.

for Non-diagonal matrices, SVD decomposition will be complex but method will same.

$\rightarrow$  Decompose any matrix into three other matrices, representing orthogonal, transformation and Scaling.



# Poornima

COLLEGE OF ENGINEERING

## DETAILED LECTURE NOTES

The diagonal elements of  $\Sigma$  are known as singular values and are non-negative. The columns of  $U$  and  $V^T$  are orthogonal orthonormal eigenvectors of  $A^T A$  and  $A A^T$ , respectively.

$$M = [U_1 \ U_2 \ \dots \ U_n] \begin{bmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_n \end{bmatrix} [U_1 \ U_2 \ \dots \ U_n]^T$$

eigenvectors                          eigenvalues

for eg:

$$\begin{bmatrix} 5 & 3 \\ 3 & 5 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$$

$$U^T = I$$

$$U U^T = I$$

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

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

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

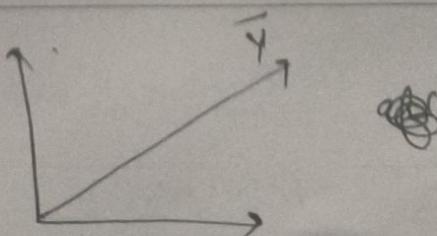
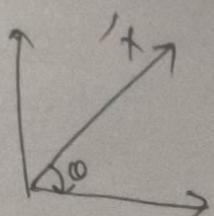
$$A \cdot A^T = U \Sigma V^T \Sigma^T U^T$$

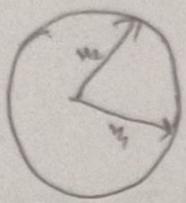
rotation . stretch . rotation

$$\begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix} \leftarrow \text{unitary transformation}$$

Stretchy

rotation





after  
stretching



$$\text{Ex:- } A = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix} \quad \begin{array}{l} \text{use complete} \\ \text{ev, } \vec{e}_1, \vec{e}_2 \end{array}$$

$$A^T = \begin{bmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{bmatrix} \quad \begin{array}{l} \text{from defn} \\ A^T \cdot A = V \end{array}$$

$$A \cdot A^T = \begin{bmatrix} 11 & 1 \\ 1 & 11 \end{bmatrix} \quad \begin{array}{l} \cancel{A^T} \\ A \cdot A^T = V \end{array}$$

Characteristic eq.

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

$$\begin{bmatrix} (11 - \lambda) & 1 \\ 1 & (11 - \lambda) \end{bmatrix} = 0$$

$$(11 - \lambda)^2 - 1^2 = 0$$

$$(11 - \lambda + 1)(11 - \lambda - 1) = 0$$

$$\begin{aligned} 11 - \lambda + 1 &= 0 \Rightarrow \boxed{\lambda_2 = 12} \\ 11 - \lambda - 1 &= 0 \Rightarrow \boxed{\lambda_1 = 10} \quad \text{ev} \end{aligned}$$

Put the  $\lambda$  value in eq.

$$\begin{aligned} \cancel{\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}} &\stackrel{\lambda_2 = 12}{=} \sqrt{1^2 + 1^2} \\ &\stackrel{\lambda_1 = 10}{=} \sqrt{1^2 + (-1)^2} \end{aligned}$$

Orthogonal base with the help of Gram-Schmidt orthogonalization process.



# POORNIMA

COLLEGE OF ENGINEERING

## DETAILED LECTURE NOTES

$$\begin{array}{c} \cancel{\sqrt{1+1}} \\ \cancel{\sqrt{1+1}} \end{array} = \cancel{1-(+1)^2}$$

PAGE NO. ....

$$\begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} = 0$$

$$A = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix}$$

$$A^T = \begin{bmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{bmatrix} \quad (A - dI) \times 20$$

$$A^T \cdot A = \begin{bmatrix} 10 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{bmatrix}$$

minor by diagonal

$$\begin{bmatrix} 10-d & 0 & 2 \\ 0 & (10-d) & 4 \\ 2 & 4 & (2+d) \end{bmatrix}$$

$$4 + 16 + 1w = \underline{\underline{120}}$$

$$\begin{bmatrix} 10 & 4 \\ 4 & 2 \end{bmatrix} / \begin{bmatrix} 10 & 4 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} 10 & 0 \\ 0 & 0 \end{bmatrix}$$

$$\boxed{\sqrt{3} - S_1d^2 + S_2d - S_3}$$

trace (A)  
diagonal

determinant of (A)

$$\begin{matrix} -2 & 0 & 2 \\ 0 & -2 & 2 \\ 2 & 2 & 4 \end{matrix}$$

B-

$$\text{trace} = 10 + 10 + 2 = 22$$

$$\begin{aligned} \det(A) &= 10(4) + 2(-20) \\ &= 40 - 40 = 0 \end{aligned}$$

we get

$$d^3 - 22d^2 + 120d = 0$$

$$d(d^2 - 22d + 120) = 0$$

$$\boxed{d_1 = 0}$$

$$d^2 - 22d + 120 = 0$$

$$\boxed{d_2 = 12}$$

$$\boxed{d_3 = 10}$$

$$\sqrt{8} \quad \checkmark$$

$$(2046)$$

$$20 - y$$

$$42 \quad 16$$

$$\frac{4+16}{20} \quad (4)$$

we get 3 values

$$\begin{bmatrix} -2 & 0 & 2 \\ 0 & -2 & 4 \\ 2 & 4 & -10 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} = 0$$

$$(A - dI)_X = 0$$

By Cramer's rule

$$\frac{u_1}{4} = \frac{-2u_2}{-8} = \frac{u_3}{4}$$

Divide by 4

$$1 \quad -2 \quad 1$$

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

$$V^T = \begin{bmatrix} 1 & 2 & 1 \\ 2 & 1 & 0 \\ 1 & 2 & -5 \end{bmatrix} \begin{matrix} \sqrt{6} \\ \sqrt{5} \\ \sqrt{30} \end{matrix}$$

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

$$16+y$$

$$=\sqrt{20} \quad \cancel{+16}$$

$$\sqrt{4} \quad \sqrt{8}$$

$$\sqrt{20} \quad \sqrt{4+16+16}$$

$$\sqrt{120}$$

$$\sqrt{8} \quad \sqrt{16} \quad \sqrt{20}$$

$$2\sqrt{2} \quad 2\sqrt{2} \quad 2\sqrt{5} \quad 2\sqrt{50}$$

$$\sqrt{2} \quad \sqrt{50}$$

$$\begin{bmatrix} \sqrt{2} & \sqrt{5} & \sqrt{10} \end{bmatrix}$$

$$-2u_1 + 2u_2 = 0$$



### DETAILED LECTURE NOTES

#### Principal Component Analysis:-

PAGE NO.....

For the given dataset in table, reduce the dimension from 2 to 1 using PCA.

Feature	$E_1$	$E_2$	$E_3$	$E_4$
$x$	4	8	13	7
$y$	11	4	5	14

No. of features = 2  
No. of samples = 4

Step 1: Calculate mean of each feature

$$\bar{x} = (4+8+13+7)/4 = 8$$

$$\bar{y} = 8.5$$

Step 2: computation of co-variance matrix,

$$C = \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) \\ \text{cov}(y,x) & \text{cov}(y,y) \end{bmatrix}$$

$$\text{cov}(x,x) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 = \frac{1}{3} (16 + 0 + 25 + 1) = 14$$

$$\text{cov}(x,y) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y}) = -11$$

$$\text{cov}(y,x) = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})(x_i - \bar{x}) = 14$$

$$\text{cov}(y,y) = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})^2 = 14$$

$$\text{cov}(y_1, y) = 23$$

$$\text{cov}(y, n) = -11$$

$$C = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

Step 3: Compute the eigenvector and eigenvalues of the covariance matrix.

i) Eigen value:

$$\det(C - dI) = 0$$

$$\begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix} - \begin{bmatrix} d & 0 \\ 0 & d \end{bmatrix} = 0$$

$$\frac{1}{29} \sqrt{5^2 - 4ac}$$

$$\begin{bmatrix} 14-d & -11 \\ -11 & 23-d \end{bmatrix} = 0$$

$$d^2 - 37d + 201 = 0$$

$$d_1 = 30, 38 \quad \left. \begin{array}{l} \\ d_2 = 6.61 \end{array} \right\} \text{Arrange in descending order}$$

Step 4: Computation of eigenvectors;  
eigenvector corresponding to  $d_1$  is a

$$u_1 = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

$$(C - d_1 I) u_1 = 0$$

$$\begin{bmatrix} 14-d_1 & -11 \\ -11 & 23-d_1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = 0 \Rightarrow \begin{aligned} (14-d_1)u_1 + (-11)u_2 &= 0 \quad \text{---(1)} \\ -11u_1 + (23-d_1)u_2 &= 0 \quad \text{---(2)} \end{aligned}$$



DETAILED LECTURE NOTES

PAGE NO.

frame ①

$$\frac{u_1}{11} = \frac{u_2}{14-1} \text{ at } u_1 = \begin{bmatrix} 11 \\ 14-1 \end{bmatrix}$$

$$u_1 = 11t$$

$$u_2 = (14-1)t$$

{ we consider largest value of  $t$  }

To find a unit eigenvector, we compute the length  $\|u_1\|$

$$\|u_1\| = \sqrt{11^2 + (14-1)^2}$$

$$\|u_1\| = 19.7369$$

Unit eigenvector

$$e_1 = \begin{bmatrix} 11/19.73 \\ (14-1)/19.73 \end{bmatrix}$$

$$= \begin{bmatrix} 11/19.73 \\ 14-30.38/19.73 \end{bmatrix}$$

$$e_1 = \begin{bmatrix} 0.5574 \\ -0.8203 \end{bmatrix}$$

(or  $\lambda = \mu$ )

$$e_2 = \begin{bmatrix} 0.8203 \\ 0.5574 \end{bmatrix}$$

Step 5 : Computation of first Principal Component

$$e_1^T \begin{bmatrix} u_{1k} - \bar{u}_1 \\ u_{2k} - \bar{u}_2 \end{bmatrix}$$

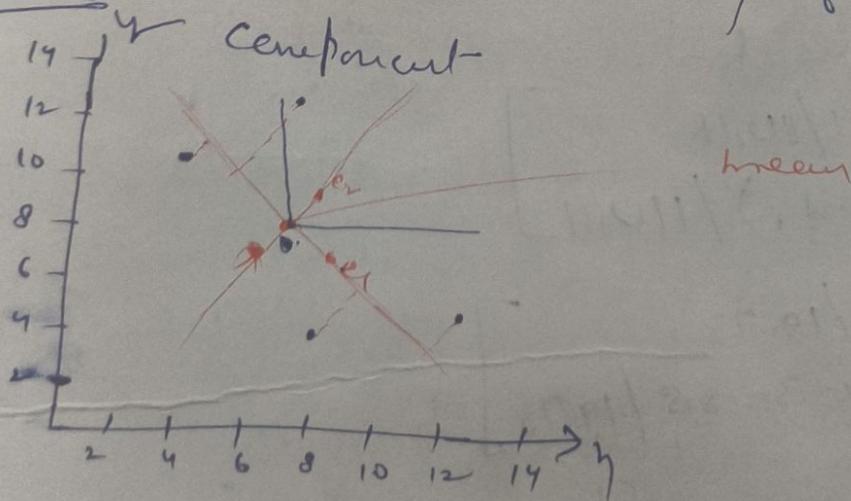
$$\begin{aligned} &= [0.5574 \quad -0.8303] \begin{bmatrix} u_{11} - \bar{u}_1 \\ u_{21} - \bar{u}_2 \end{bmatrix} \\ &= 0.5574(u_{11} - \bar{u}_1) - 0.8303(u_{21} - \bar{u}_2) \\ &= .55(4-8) - .83(11-8.5) \\ &= -4.3053 \end{aligned}$$

So

The result of calculation is summarized in table:

$u$	4	8	13	7
$y$	11	4	5	14
$P_C$	-4.3053	3.7361	5.6928	-5.1230

Step 6 : Geometric meaning of first principal component



## PCA in machine learning:

--PCA stands for Principal Component Analysis.

--PCA or Principal Component Analysis, is a common technique used in machine learning for dimensionality reduction.

--Dimensionality reduction refers to the process of reducing the number of features or variables in a dataset while preserving the most important information. This is important because datasets with a large number of features can be difficult to work with and can lead to overfitting and decreased performance of machine learning models.

--remove the problem of curse of dimensionality. The curse of dimensionality refers to the difficulties that arise when working with high-dimensional data in machine learning. As the number of

dimensions (or features) of a dataset increases, the amount of data required to obtain statistically significant results increases exponentially. This can lead to problems such as overfitting, increased computational complexity, and reduced model performance. This makes it difficult to find meaningful patterns or relationships within the data, and it becomes more challenging to distinguish between relevant and irrelevant features.

--PCA works by identifying the directions in the data that contain the most variance, and then projecting the data onto a new coordinate system that is aligned with these directions. The new coordinate system consists of a set of orthogonal axes (principal components) that capture the largest possible variance in the data. The first principal component captures the most variance, followed by the second, and so on.

Advantages of PCA Algorithm: PCA (Principal Component Analysis) is a widely used technique in machine learning for dimensionality reduction. Here are some advantages of using PCA:

1. Reduces the number of features: PCA can help reduce the number of features in a dataset while retaining most of the information. This can be particularly useful when dealing with high-dimensional data, where the number of features is much larger than the number of samples.
- 2 Improves model performance: By reducing the number of features, PCA can help improve the performance of machine learning models. This is because it reduces the risk of overfitting, which can occur when there are too many features relative to the number of samples.

3. Speeds up training: With fewer features, training machine learning models can be faster and more efficient. This is because there are fewer calculations to perform, and the resulting models are simpler and easier to interpret.

4. Helps with visualization: PCA can be used to visualize high-dimensional data in two or three dimensions. This can help identify patterns and relationships in the data that may not be apparent in higher dimensions.

5. Removes correlated features: PCA can help remove correlated features, which can cause problems for some machine learning algorithms. By removing these features, PCA can help improve the stability and accuracy of the resulting models.

data. This can lead to suboptimal results if the dataset contains outliers.

to understand the underlying patterns in the data or explain the results of a machine learning model to others.

3. Assumes linearity: PCA assumes that the underlying relationships between the features are linear. However, in some cases, the relationships may be non-linear, which can lead to suboptimal results.

4. Computationally expensive: PCA involves computing the eigenvectors and eigenvalues of the covariance matrix of the data, which can be computationally expensive for large datasets. This can make PCA impractical for some applications.

5. Sensitivity to outliers: PCA is sensitive to outliers, as outliers can have a large influence on the covariance matrix of the

Disadvantages of PCA Algorithm: While PCA (Principal Component Analysis) is a powerful technique for dimensionality reduction in machine learning, it also has some disadvantages. Here are some of the main disadvantages of PCA:

1. Information loss: PCA works by finding a new set of features that capture most of the variance in the data. However, this can result in some information loss, as the new features may not capture all of the information in the original features. This can be particularly problematic if the discarded information is important for the task at hand.

2. Interpretability: The new features generated by PCA are linear combinations of the original features, which can make them difficult to interpret. This can be a problem if we want