

Data Analysis - Process of examining raw data with the purpose of drawing conclusion about that information.

Vector Space: also known as linear Spac., is a fundamental concept in linear Algebra.

It is a collection of vectors, which are objects that can be added together and multiplied ("scaled") by numbers, called scalars, in this context.

- A vector space must satisfy a few key properties, related to vector addition & scalar multiplication.

1. Closure under Addition

If v_1 & v_2 are vectors in space, then their sum $v_1 + v_2$ is also in space.
e.g. $v_1 = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$ $v_1 + v_2 = \begin{bmatrix} 3+(-1) \\ 4+2 \end{bmatrix} = \begin{bmatrix} 2 \\ 6 \end{bmatrix}$
 $v_2 = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$

2. Closure under Scalar Multiplication

- when you multiply any vector (or matrix) within a set by a scalar (a real no.), the resulting vector (or matrix) also remains within a set.
e.g. if v is a vector and c is a scalar, then the product $c \cdot v$ is also in the space.

→ operation on data structures such as vectors, matrices, or datasets)

- Specifically, a dataset or a vector space is closed under scalar multiplication, if scaling the data by a constant results in a valid data point or a data structure, i.e. still consistent with the structure of original datasets.

e.g. Feature Scaling in ML

Consider a scenario, where you are working with a dataset of features for a ML model. Let's say each feature is represented by a vector in a vector space. The vector space is closed under scalar multiplication, if scaling any feature (multiplying it by a scalar) still results in a valid feature vector.

Age	Income
95	50000
30	60000
22	45000

You have a dataset with two features Age & Income of a group of individuals. This can be represented as a Matrix:

$$X = \begin{pmatrix} 95 & 50000 \\ 30 & 60000 \\ 22 & 45000 \end{pmatrix}$$

Now let's say you scale the Age feature by a constant scalar value $c=2$. The new feature matrix becomes

$$X_{\text{scaled}} = \begin{pmatrix} 50 & 50000 \\ 60 & 60000 \\ 44 & 45000 \end{pmatrix}$$

In this case, you have scaled the Age column, but Income remains unchanged.

The new Matrix is still a valid feature matrix for a ML model, as Dataset is closed under scalar multiplication, so multiplying by a scalar still results in valid data.

↳ Case where Scale Multi/Divide is not valid
e.g. Dealing with Binary variable (e.g. Gender, where 0 = female and 1 = male)

$$\text{Consider a dataset: } \begin{matrix} \text{Gender} \\ 0 \\ 1 \end{matrix}$$

If we multiply the Gender feature by a scalar $c=2$, we get

$$X_{\text{scaled}} = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$$

The new values of Gender (0, 2, 0) are no longer valid binary values.

This means the dataset is not closed under scalar multiplication in this case.

Scalar multiplication destroys the integrity of a data, because it introduces invalid values that fall outside the expected range of the variable.

⇒ In Data Analysis - Closure under scalar multiplication ensures that transformation like scaling preserve the validity & consistency of dataset.
It depends on the nature of a data and transformation applied.

3. Commutativity Addition - in Data Analytics
 when you add elements (matrices or vectors) doesn't affect the result. e.g. $a+b = b+a$
 This means that the sum of two elements is independent of the order in which they are added. This property is crucial when manipulating data in datasets, as it ensures consistency like aggregating or combining data.

Ex: working with Vectors (Feature Data)

Consider two data points (vectors) representing features of two individuals in a dataset, say age & income.

• vector 1 (Individual 1): $v_1 = \begin{pmatrix} 25 \\ 50000 \end{pmatrix}$

• vector 2 (Individual 2): $v_2 = \begin{pmatrix} 30 \\ 60000 \end{pmatrix}$

If we add these vectors, we get:
 $v_1 + v_2 = \begin{pmatrix} 25 \\ 50000 \end{pmatrix} + \begin{pmatrix} 30 \\ 60000 \end{pmatrix} = \begin{pmatrix} 55 \\ 110000 \end{pmatrix}$

Now, if we reverse order of addition
 $v_2 + v_1 = \begin{pmatrix} 30 \\ 60000 \end{pmatrix} + \begin{pmatrix} 25 \\ 50000 \end{pmatrix} = \begin{pmatrix} 55 \\ 110000 \end{pmatrix}$

→ As we can see, the order of addition doesn't affect the result.
 ∴ the operation of addition is commutative.

⇒ Real life Application: Aggregation of Data.
 Imagine you are analyzing sales data across multiple regions or stores. You may need to aggregate sales data from different stores. The commutative property of addition ensures that it doesn't matter which store's data you add first, total sales remain same.

• Sales 1: $v_1 = [500, 400, 600]$ → Aggregating Sales Data
 • Sales 2: $v_2 = [300, 200, 150]$ $v_1 + v_2 = [800, 600, 750]$

And, if we reverse the order
 $v_2 + v_1 = [800, 600, 750] \Rightarrow$ Again, the order of addition doesn't matter, showing that commutativity holds when aggregating data.

why is commutativity so important?

- Consistency - When performing operation like summation, the commutative property ensures that the result will be consistent, regardless of order of operation. This is important when aggregating data or combining features.
- Efficiency in Calculation - In large datasets, algorithms can take advantage of commutativity to optimize calculation, e.g. computation across multiple processors, distributing the computation across multiple processes.
- Data Integrity - Commutativity guarantees that when you are manipulating data (e.g., sum, total), the order in which data is processed does not affect the outcome.

4. Associativity of Addition:
It states that when you add three or more numbers, the grouping of those numbers does not affect the result.

$(a+b)+c = a+(b+c)$
When performing operation like data aggregation, feature engineering or summarizing multiple data points.

Cg. Working with vectors (Data Points):

Consider three data points (represented as vectors), each corresponding to a person's age and income.

$$v_1 = \begin{pmatrix} 95 \\ 10000 \end{pmatrix}; v_2 = \begin{pmatrix} 30 \\ 60000 \end{pmatrix}; v_3 = \begin{pmatrix} 22 \\ 45000 \end{pmatrix}$$

FIRST GROUPING
 $(v_1 + v_2) + v_3 = \left(\begin{pmatrix} 95 \\ 10000 \end{pmatrix} + \begin{pmatrix} 30 \\ 60000 \end{pmatrix} \right) + \begin{pmatrix} 22 \\ 45000 \end{pmatrix} = \frac{77}{155000}$

SECOND GROUPING
 $v_1 + (v_2 + v_3) = \begin{pmatrix} 95 \\ 10000 \end{pmatrix} + \left(\begin{pmatrix} 30 \\ 60000 \end{pmatrix} + \begin{pmatrix} 22 \\ 45000 \end{pmatrix} \right) =$

But, Addition is Associative, and the order in which we group the vectors for addition does not change the final results.

Working with Dataframes - First + Grouping:
Working with Dataframes - First + Grouping.
 $(df_1 + df_2) + df_3 = \left(\begin{pmatrix} 200 & 250 & 300 \\ 150 & 200 & 250 \end{pmatrix} + \begin{pmatrix} 300 & 350 & 400 \\ 200 & 250 & 300 \end{pmatrix} \right) + \begin{pmatrix} 100 & 150 & 200 \\ 50 & 100 & 150 \end{pmatrix}$
 $= \begin{pmatrix} 500 & 600 & 700 \\ 350 & 450 & 550 \end{pmatrix} + \begin{pmatrix} 100 & 150 & 200 \\ 50 & 100 & 150 \end{pmatrix} = \begin{pmatrix} 600 & 750 & 900 \\ 400 & 550 & 700 \end{pmatrix}$

Second Grouping
 $df_1 + (df_2 + df_3) = \checkmark$

\Rightarrow In Data Analytics, the Associativity of Addition allows flexibility and consistency when performing operation such as Aggregating, combining or summing data points, vector or matrices.

5. Existence of Additive Inverse

Additive inverse of a number is the value that, when added to the original number, result in zero.

Every number has an additive inverse, and it is found by simply changing the sign of the original number.

$$a + (-a) = 0$$

Suppose we have 2×2 matrix representing sales data for two food items across two stores.

$$M = \begin{pmatrix} 100 & 150 \\ 200 & 250 \end{pmatrix}$$

The additive inverse of Matrix M is simply the matrix with all values negated

$$-M = \begin{pmatrix} -100 & -150 \\ -200 & -250 \end{pmatrix}$$

$$\text{Now let's add } M \text{ and } -M = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$$

SubSpace :

In a high dimensional dataset, many dimensions might be irrelevant or noisy, masking underlying patterns and clusters.

Subspace Analysis, focuses on finding meaningful clusters within specific subsets (Subspace) of the original dimension.

↳ Why use Subspace Analysis?

- It helps to reduce the complexity of high dimensional data by focusing on the most relevant features.
- Subspace clustering algorithms can identify clusters that might be hidden or obscured in the full dataset.
- Subspace methods can also be used for outliers detection by identifying data points that deviate significantly from the subspace characteristics.

Eg. Customer Segmentation in a E-commerce Platform.
[How subspaces are used to simplify complex, high-dimensional data, uncovers patterns & improve decision making].

Suppose you are working for an e-commerce company that sells a wide variety of products online with various features like.

You have customer data that includes features like Total Amount Spent, No. of Purchases, Total Amount Spent, Time Spent on Site, Age, Gender, Product Category Purchased, Geographic region (City, State, Country), Product Category Purchased, etc.

This dataset has many dimensions (features) and it would be hard to manually analyze and understand customers behavior directly from this data.

However, we can use subspaces to reduce data dimensionality and identify key patterns and create meaningful segments of customers.

5. Initial Dataset: High Dimensional Feature Space

<u>CustId</u>	<u>Age</u>	<u>Gender</u>	<u>Time spent</u>	<u>Purchase</u>	<u>Total AmntSpent</u>	<u>Product Category</u>	<u>Region</u>
1.	25	M	45	5	150	Electronics	DefD, Cust2D
2.	34	F	30	3	90	HomeGard	Mumbai
3.	45	M	60	10	350	Sport	DefD
4.	27	F	20	2	50	Apparel	Chennai
5.	36	M	50	7	220	Sport	DefD ..

2. Challenges: High dimensionality.
 While this data can tell you some information about individual customers, it's hard to find patterns.
 Instead of analyzing the full 7-features (dimensional space), we can reduce the complexity by finding a subspace that better capture key features of customer behavior.

3. Dimensionality Reduction

- PCA (Principal Component Analysis)
 - ↳ a technique used to identify directions (or components) in which the data varies the most and reduce the dimensionality of data.
 - To perform PCA, on the given dataset, we need to focus on numerical columns that PCA can process.
 - PCA works on Numerical data, so we can exclude categorical data such as Gender, Product Category Purchase and Region.
- Steps for PCA

1. Prepare the data - organize numerical column for PCA and standardize the data.

2. Standardize the data - Standardization is important, when features have different scales, so we normalize the data to have a mean '0' and a standard deviation of 1.

3. Compute Covariance Matrix & Eigenvalues & Eigen Vectors. The eigenvectors corresponding to the largest eigenvalues are the principal components.

4. Find Eigenvalues & Eigen Vectors. The eigenvectors corresponding to the largest eigenvalues are the principal components.

5. Sort Eigenvalues & Eigen Vectors. The eigenvectors corresponding to the largest eigenvalues are the principal components.

6. Project the Data - Transform the data onto the new set of principal components.

P(A)) Step 1 : Standardization - calculate Mean & std. deviation of each column (Age, Time Spent, Purchase, Total Amnt spent)

CustID	Age	Time Spent	Purchase	Total Amnt Spent	
1	25	45	5	150	Simplify Problem
2	34	30	3	90	
3	45	60	10	350	
4	27	90	2	50	
5	36	50	7	220	

- Mean - is the sum of value divided by the number of values.
It represents central tendency of data.
- Mean = $\frac{\sum x_i}{N}$
 - x_i - each data point
 - N - No. of Data Point (In this case, 5)

- Standard Deviation - tell how you spread out the values are from the Mean. A high standard deviation mean the data points are spread out, while a low standard deviation means they are close to mean.

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2}$$

- x_i - each data point
- \bar{x} - Mean of data point
- N - No. of data point (here, 5)

1. Age : Mean = $\frac{25 + 34 + 45 + 27 + 36}{5} = 33.4$

- Standard deviation - First, calculate the squared difference from Mean
 $(25 - 33.4)^2 = 70.56$; $(34 - 33.4)^2 = 0.36$; $(45 - 33.4)^2 = 133.96$.
 $(27 - 33.4)^2 = 40.96$; $(36 - 33.4)^2 = 6.76$
- Now divide by $N = 5$, and take the square root.
 Sum of squared difference = $70.56 + 0.36 + 133.96 + 40.96 + 6.76 = 252.6$

$$\text{standard deviation} = \sqrt{\frac{252.6}{5}} = \sqrt{50.52} = 7.12.$$

2. Time Spent : value: 45, 30, 60, 20, 50
 $\text{Mean} = \frac{45+30+60+20+50}{5} = \frac{205}{5} = 41$

Standard deviation :
 $(45-41)^2 = 16$, $(30-41)^2 = 121$, $(60-41)^2 = 361$

$(20-41)^2 = 121$, $(50-41)^2 = 81$

Sum of squared difference :

$$16 + 121 + 361 + 121 + 81 = \underline{\underline{1020}}$$

$$\text{std. Deviation} = \sqrt{\frac{1020}{5}} = \sqrt{204} = \boxed{14.28}$$

3. Purchase : value: 5, 3, 10, 2, 7.

$$\text{Mean} = \frac{5+3+10+2+7}{5} = \frac{27}{5} = 5.4$$

Sum of squared differences - First calculate sum of squared differences.

$$(5-5.4)^2 = 0.16, (3-5.4)^2 = \underline{\underline{5.76}}, (10-5.4)^2 = \underline{\underline{21.16}}$$

$$(2-5.4)^2 = \underline{\underline{11.56}}, (7-5.4)^2 = \underline{\underline{2.56}}$$

Sum of squared differences

$$0.16 + 5.76 + 21.16 + 11.56 + 2.56 = \underline{\underline{41.9}}$$

Now divide by $N=5$ and take the square root

$$\text{std. Deviation} = \sqrt{\frac{41.9}{5}} = \sqrt{8.38} = \boxed{2.87}$$

big
data
set

4. Total Amt. Spent: Value: 150, 90, 350, 50, 20
 $\text{mean} = \frac{150 + 90 + 350 + 50 + 20}{5} = \frac{860}{5} = 172$

Std. deviation - First calculate squared difference from Mean,
 $(150 - 172)^2 = 484$, $(90 - 172)^2 = 6724$, $(350 - 172)^2 = 31584$
 $(50 - 172)^2 = 14884$, $(20 - 172)^2 = 2304$

Sum of squared differences:
 $484 + 6724 + 31584 + 14884 + 2304 = 55176$
 Now divide by $N=5$, and take the square Root:
 $\text{std.deviation} = \sqrt{\frac{55176}{5}} = \sqrt{11035.2} = 105.06$

⇒ Summary

<u>Column</u>	<u>Mean</u>	<u>Standard Deviation</u>
Age	33.4	7.12
TimeSpent	41	14.28
Purchase	5.4	2.87
Total Amt Spent	172	105.06

Step 1:
 Mean and std.
 deviation

Step 2: Covariance Matrix - After Standardization, calculate the Covariance Matrix to see, how each feature correlate with other.

Now, we need to calculate the covariance between each pair of variables in the data set. Covariance Matrix will have following form:

$$\begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) & \text{cov}(x,z) & \text{cov}(x,w) \\ \text{cov}(y,x) & \text{cov}(y,y) & \text{cov}(y,z) & \text{cov}(y,w) \\ \text{cov}(z,x) & \text{cov}(z,y) & \text{cov}(z,z) & \text{cov}(z,w) \\ \text{cov}(w,x) & \text{cov}(w,y) & \text{cov}(w,z) & \text{cov}(w,w) \end{bmatrix}$$

- $x = \text{Age}$
- $y = \text{TimeSpent}$
- $z = \text{Purchase}$
- $w = \text{Total Amt. Spent}$

The Covariance between two variables X and Y is given by

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

- x_i & y_i are individual data points of each variable.
- \bar{x} and \bar{y} are the mean of the variables.
- n is the no. of data points.

Variance of a variable is just the covariance of the variable with itself:

$$\text{Cov}(X, X) = \text{Var}(X)$$

which is equivalent to the square of the standard deviation of the variable:

\text{Var}(X) = \sigma^2_x

So, if we assume zero correlation, the covariance matrix would be:

$$\begin{bmatrix} \sigma^2_{\text{Age}} & 0 & 0 & 0 \\ 0 & \sigma^2_{\text{Time Spent}} & 0 & 0 \\ 0 & 0 & \sigma^2_{\text{Purchase}} & 0 \\ 0 & 0 & 0 & \sigma^2_{\text{Total Time Spent}} \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} (7.12)^2 & 0 & 0 & 0 \\ 0 & (14.28)^2 & 0 & 0 \\ 0 & 0 & (2.87)^2 & 0 \\ 0 & 0 & 0 & (105.06)^2 \end{bmatrix}$$

This, Covariance Matrix is:

$$\begin{bmatrix} 50.6944 & 0 & 0 & 0 \\ 0 & 204.9184 & 0 & 0 \\ 0 & 0 & 8.2364 & 0 \\ 0 & 0 & 0 & 11043.0036 \end{bmatrix}$$

8

Interpretation of Covariance Matrix

- The diagonal elements represents the variance of each variable, showing how spread out the data for each variable is.
 - The off-diagonal elements (which are all zero) represent the covariance b/w the variables.
- Since we assumed zero correlation, the covariance value between different variables is zero.

If correlation exist - if you had raw data or correlation coefficients between the variables, you could calculate the off-diagonal elements by multiplying std. deviations of the respective variable and the correlation coefficient.

e.g., if the correlation between Age and time spent 0.3, the covariance would be:

$$\text{Cov}(Age, Time Spent) = \text{Correlation}(Age, Time Spent) \times \sigma_{Age} \times \sigma_{Time Spent}$$

Step 3: Eigenvalues and Eigenvectors

[Computst. of Covariance Matrix]

$$\begin{bmatrix} 50.6944 & 0 & 0 & 0 \\ 0 & 404.9184 & 0 & 0 \\ 0 & 0 & 8.2364 & 0 \\ 0 & 0 & 0 & 11043.0036 \end{bmatrix}$$

This is a diagonal matrix, where all off-diagonal elements are zero, indicating no correlation between the variables.

Eigenvalue - each eigenvalue corresponds to the variance along a principal component direction.

Eigenvector - Since Covariance Matrix is diagonal, the eigenvectors will be the standard basis vector.

Calculate the Eigenvalues

For the diagonal Matrix, Eigenvalue are simply diagonal values themselves.

$$\lambda_1 = 50.6944, \lambda_2 = 204.9184, \lambda_3 = 8.269, \lambda_4 = 110.43.$$

Calculate the Eigenvectors - For a diagonal matrix, the Eigenvectors are straightforward. Each Eigen vector corresponds to one of the dimensions in data space (Feature: Age, TimeSpent, Purchase, Total Amt Spent).

$$v_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, v_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, v_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, v_4 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

- Each Eigen vector corresponds to the direction of one feature.
- v_1 - corresponds to Age; v_2 - corresponds to TimeSpent.
- v_3 - corresponds to Purchase; v_4 - corresponds to Total Amt Spent.

Eigenvalues

$\lambda_1 = 50.6944$, indicates the variance explained by the first Principal Component (Age)

- $\lambda_2 = 204.9184$, — 2nd, principal component (TimeSpent)
- $\lambda_3 = 8.269$, — 3rd Principal Component (Purchase)
- $\lambda_4 = 110.43.0036$, — 4th, Principal Component (Total Amt Spent)

Eigenvector - being standard basis vectors, indicate that the features (Age, TimeSpent, Purchases, Total Amt Spent) are independent of each other in this case. There is no mixing between them along the principal components.

$$v_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} (\text{Age}), v_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} (\text{TimeSpent}), v_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} (\text{Purchase}), v_4 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} (\text{Total Amt. Spent})$$

These results show that each feature in the dataset contributes to a distinct principal component, with no mixing or interaction between them.

Independence of Vector

Vector do not "depend" on each other for their existence, and none of them can be written as a combination of others. e.g. Let's consider a dataset with two features (e.g. Age & Income of individuals) represented as two vector.

$$\text{vector 1 (Age)} \quad v_1 = \begin{pmatrix} 25 \\ 30 \\ 22 \\ 28 \end{pmatrix}$$

$$\text{vector 2 (Income)} \quad v_2 = \begin{pmatrix} 50000 \\ 60000 \\ 45000 \\ 55000 \end{pmatrix}$$

Now, let's see whether these two vectors are linearly independent.

To check if they are linearly independent, we could check if there exists a nontrivial solution to the equation:

$$c_1 v_1 + c_2 v_2 = 0$$

This gives us the following system of equations:

$$c_1 \begin{pmatrix} 25 \\ 30 \\ 22 \\ 28 \end{pmatrix} + c_2 \begin{pmatrix} 50000 \\ 60000 \\ 45000 \\ 55000 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

which translates to

$$\begin{pmatrix} 25c_1 + 50000c_2 \\ 30c_1 + 60000c_2 \\ 22c_1 + 45000c_2 \\ 28c_1 + 55000c_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

For these vectors to be linearly independent, the only solution to this system of equations should be $c_1=0$ and $c_2=0$. If we attempt to solve this system, we will find that there is no nontrivial solution whose both c_1 and c_2 are nonzero. Thus, these vectors are linearly independent, meaning that the age and income features are not directly related to each other in such a way that one can be written as a multiple of other. This implies that age & income are distinct features in this dataset, each contributing unique information to the analysis.

\Rightarrow Basis & dimension allow you to understand the structure of your data, reduce complexity and analysis & modeling efficiency.

- Basis - A ~~tiny~~ minimal set of features that can describe your data.
- Dimension - The no. of features needed to describe your data. (i.e. - Basis of the space)

\Rightarrow Techniques like PCA and feature Selection help reduce the dimension by finding a new basis with fewer features.

e.g. in context of data Analytics

- Suppose you have a dataset with multiple features (e.g. Age, Income, education). These features can be thought of as vectors in a higher dimensional space.
- A basis would be a set of features that can describe all possible data points in your dataset.
You might not need all features, and some features could be redundant, meaning they are not part of a "minimal" basis.

e.g. Dataset with three variables

- x_1 : Age ; x_2 : Income ; x_3 : Edu. Level
— if you find that x_1 & x_2 are enough to describe all data points (say, education level just a function of age & income), then x_1 & x_2 form a basis for the data space, and x_3 would be redundant.

\Rightarrow if your dataset contains 5 features (Age, income, education, years of experience, and city), then your vector space has a dimension of 5.

Data Analysis

L process of examining raw data and the purpose of drawing conclusions about that information.