# *Shallow Copy vs Deep Copy

↓

→ pointer points to the Same Copy of objects Of class

→ Store references of object to original Memory Address

→ reflect changes made to the new (copied object cf Originally object

→ faster

↘

Create copy of each object inside of Class

→ Store Copies of object's value.
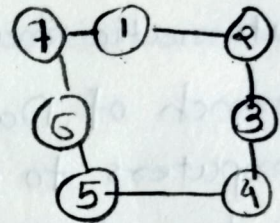
→ Doesn't reflect Changes

→ comparitively slower

---

ndarray vs array

# Data Science Life Cycle :—

- Refer Series <u>of steps</u> — follow by data Scientists — working of Data Driven program
- <u>Include</u> = data — cleaning, preparation, modelling, model, etc.
- lengthy procedure & quitly take few months.
- essential to have generic structure — CRISP — DM Framework.

1. Business Understanding
2. Data Understanding
3. Data Preparation
4. Explovitory Data Analysis
5. Data Modeling
6. Data Evulation
7. Model Deployment

## 1. Business Understanding
- Enterprise goal
- aim of analysis
- aim of Evaluation
→ Prickction

## 2. Data Understanding.
- Series of all reachable data
- Explore information using graphical plot
- Exploring data

## 3. preparation of Data.
- choosing applicable Data, Integrate the data
→ formate data into preferred Structure
→ Time Consuming, arguably the most esential Step
→ Model will be accurate as your data

## 4. Explorcatory Data Analysis.
- getting some concept about the answer and element affecting it.
- Discover each and Every char indi. & by means Combine them with diff. feature.

## 6. Data M Evalution.
- gecised up to deploy

## 5. Data Modeling :—
- Company heart of data Analysis
- Organized Data. input |pref. output
- Selecting Model
- hyperparameter

## 7. Model Deployment

# → Application of Datascience.

- ~~Search Engine~~
- → gaming
- → Health Sector

- finance
- Business
- Airline

- Medical
- Image Recognize
- E-commerce

# → NLP

- Natural Language Processing
- Automatic manipulation of Natural languages
- Branch of Data Science that focus on training computers to process and interpret conversation in text format in a way human do by listen.
- → filling gap between Data Science & human language

→ Difficult & challenging during Development as computer require humans to interect with them Using programming languages like Java, python, etc.

↓

Structure & unambiguous

⇒ Use Case

1. spam detection
2. Machine translation
3. Virtual Agents & chatbot
4. Social Media Sentiment Analysis
5. Text Summerization

# ⭐ Computer Vision

→ field of artificial Intelligence
  - train Computers to intrup & understand
    - visual world

→ Machine - accuratly idetify objects - then react
  ___ Work ___
→ need lots of data
→ performed over & over analyses of data
Conclusion→ - until it recognozj output

  Ex & tire

→ 2 Essential technologies :- Machine Learning (Deeplea.)
                           :- Convolution neural Network
                              (CNN)

## Application

① IBM - 2018, My momouts - Masters golf & turnament

② google translator

③ Self- Driving vehicals

④ IBM- partnership - Verizon to bring AI to Edge.

# Big Data

→ Collection of huge data, yet growing rapidly with time.

→ Complex & huge in size — none Data Managment System can store it.

→ System - process and store huge data become common component in Big Data Management in organization - Combined tool that support Big Data uses.

→ Charachatrize by three V's :-

① large Volume - memy Enviroment
② Coide Varitety - Stored in Big Data System
③ Velocity - much data is generated, collected & processed.

⟹ Application

→ Banking        → agriculture      → finance
→ Education      → Media            → E-commerce
→ Healthcare     → Goverment        → Retail
→ Agri

⟹ Issue In Big Day Data Science.

- empowring compnies
- data constantly - has to be handle & cant ignored
- Data Scientist Collect Data Set, remove Data, Analyse.

① Identifying problem
② Fiding appropirte data
③ Coorkforce
④ Clensing.

# Scalar

- Singal Numerical Value

  Ex $[0]$ $[1]$

# Vector

- Array of Numbers

  Ex $[0,1,2]$ $[3,8,9]$

# Matrix

- 2D Array of Shape with $m$ rows & $n$ columns.

  Ex $\begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix}$ $\begin{bmatrix} 9 & 0 \\ 10 & 1 \end{bmatrix}$

# Descriptiv Statistics

- Describing & Summerizing data numerically.

- 2 approches

  ① Quentitve Approche - numerically
  ② Visual Approch - illustrate data with graph, charts, etc.

# Types of Mesure

① Central Data - Cente of data
   (tendency)
   - Mean, Median, Mode

② Variblity - Spread of data
   - Variance & Standard deviation

③ Correlation (or)
   Joint Valiblity :- relation between pair of veriables.
   - correlation coefficient , covariation

# Mean

- The _average_
- Central value of _finite set of number_

$$\mu = \frac{\sum_{i=1}^{N} x_i}{N}$$

Example: $x_1, x_2, \dots x_N$

## Median

- middle element of dataset.
- Stored in incresing or decreasing
- → odd. ⇒ middle value.
- → EVEN ⇒ two value at midde position

## Mode

→ Most frequently

→ isn't such single value, then set is multimodel Since it has multiple model value.

## Variance

→ Messure how far the data points are spread out from the average value, and is equal to the sum of squares of diff. between the data value and the average.

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

## Standard deviation

⇒ Squre root of varience and measure the extent to which data varies from its Mean

→ prefered over varience.

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

$\longrightarrow$ popular tools used in data Science.

- Data-preprocessing & Analysis
$\rightarrow$ Data Exploration & Visualization
$\rightarrow$ Parallel and distributed Computing incase of BigData

$\longrightarrow$ Python as programming language.
$\rightarrow$ Support Multiple programming paradigm
$\rightarrow$ Dynamic typing
$\rightarrow$ Reference Counts
$\rightarrow$ Late biding

* 20 Algorithm as described in zen of python by Tim Peters.

$\longrightarrow$ Data Science Application
1. Search Engine
2. Transport
3. Finance
4. E-commerce
5. Healthcare
6. Image recognization

---

## Covariance

$\longrightarrow$ Mesure of the joint variability of two random variables & Describes relationship between them

$\rightarrow$ Defined as expected value of the product of the two random variables's deviation from their means.

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

number of Data → Data Value of x → Mean of x

① positive
$\rightarrow$ heading in Same direction

② Negative
$\rightarrow$ heading in opposite

# EDA [Exploratory Data Analysis]

→ Critical Step - investigate and summarize - main char. of Data Set

→ Help to understand data, identify pattern & uncover insights

## 1. Loading the Data

- Start by Loading Dataset.

```
import pandas as pd
df = pd.read-csv ('your-dataset.csv')
```

## 2. Inspecting the Data :-

- Check first few row and basic info

```
print (df.head())
print (df.info())
```

## 3. Handling Missing Value.

- Identify & Handle Missing Value.

```
print(df.isnull().sum())
df.fillna (df.mean(), inplace = True)
```

## 4. Descriptive Statistics :

- calculate summary Statics to understand data's central tendency & varibility

```
print (df.describe())
```

## 5. Data Visualization :

- Create visualization to better understand the data.
- lib :- Matplotlib , Seaborn

```
import matplotlib.pyplot as plt
import Seaborn as sns

plt.hist(df ['Age'])
plt.xlable ('Age')
plt.ylable ('Frequency')
plt.show()
```

## 6. Feature Relationships.

- Explore relation between variable.

```
correlation-matrix = df.corr()
sns.heatmap (correlation-matrix, annot=True)
plt.show()
```

## 7. Data transformation

- perform Data transform is needed

```
df= pd.get-dummies (df, columns=['Category'])
from sklearn.preprocessing import MinMax Scaler
Scaler = MinMaxScaler()
df['Age'] = Scaler.fit-transform (df[['Age']])
```

## 8. Outerliner Detection:

- Identify + Deal with outliner in data

## 9. feature Engineering:

- Create new feature or drive insight from existing features

```
df ['Total-score'] = df['Score1'] +df['score2']
```

## 10. final Summary:-

Provide a final Summary of your finding, Insight, and any recommendation based on your analysis.