# Exploratory Data Analysis (EDA):

By: Loga Aswin

It refers to the method of studying and exploring records to apprehend their predominant traits, discover patterns, locate outliers & identify relationships between variables.

- ## Goals of EDA:

1. **Data cleaning**: Handle missing values, duplicates, outliers & handle Categorical Values.

2. **Data Visualization**: Visual techniques represent Statistics Graphically. Histogram, box plots, Scatter plot, line plot, heat map & bar charts to identify Styles, trends & relationships with facts.

3. **Feature Engineering**: Contain Scaling, Normalization, binning, encoding Variables.

4. **Correlation & Relationships**: Allow discover relationship & dependencies between Variables. Correlation analysis, Scatter plot & pass tabulation.

5. **Data Segmentation**: Divide information into Significant Segments based totally on Swe-Standard or traits.

6. **Hypothesis Generation**: Generating hypothesis / Studies Questions based on preliminary exploration of data.

7. **Data Quality Assessment**: Permits assessing nice & reliability of the info.
   ↳ Involve checking record integrity, consistency of Accuracy to make into Suitable for Analysis.

① .

# ⊙ Handling Missing Values :

- **isnull ()** - check any missing values in dataset.

  → df. isnull . sum ()

  O/P → Return sum of np.nan (NULL) values in each column.

- **fillna ()** - fill value at NULL Places.

  → df [' column']. fillna (value, inplace = True)

- **replace ()** - used to replace values in dataset.

  → df ['col']. replace (np.nan , values)

- **dropna ()** - Drop records with NULL Values.

  → df ['col']. dropna (axis = 0, how = 'any')

- **duplicated ()** - Checks if duplicates Present in dataset.

  → df. duplicated ()

  O/P : Returns total no. of duplicate rows.

- **drop. duplicates** - Drop the Duplicate rows.

  → df. drop.duplicates (Keep = 'first', inplace = True)

  O/P : keep first copy & remove all other duplicates.

⊙ **Data Encoding :**

→ Encode Categorical data into numerical values .

\* One - hot Encoding / Label Encoding :

     ↝ from sklearn.preprocessing import Label Encoder

     ↝ encoder = Label Encoder ()

     ↝ df ['col'] = encoder. fit - transform (df ['col'])

O/p → Assign num to each Category starting from 0 .


⊙ **Data Visualization :**

↳ Analyze data in the form of Graphs / maps, easy understand trends / patterns .

① . **Boxplot** —

         sns . boxplot (x = 'col1' , y = 'col 2' , data =df)

② . **Pairplot** () - pairwise distribution in dataset .

         sns. pairplot (df , hue = 'col', height = 2)

③ . **Histogram** : Count the numeric values in a dataset .

         sns. histplot (x = 'col', data = df)