

EDA & FE

available

EDA \rightarrow Exploratory Data Analysis
FE \rightarrow Feature engineering

DS life cycle

1] Data collection or ingestion



EDA [Exploratory Data Analysis]



Pre-processing the data



Model building



Evaluate the Model

Statistics

* Statistics is science of collecting the data & organize, Analyze the data.

* After this we will get some insight from the data after the statistics implementation.

So, the first point is statistics Analysis in EDA section.

Real time Scenario why we need ML

Example

Sale of Product

↳ Sale is going down



Reasons

Product, paying to customer,
leadership, marketing, competitors



going to collect the dataset
to determine the solution.

with help of Analysis.

Role of persons (Involved in any kind of
above problem)

1) Project Manager

2) Business analyst

3) Data Scientist

↳ give the
conclusion.

Any kind of data set we have to
do analysis

Famous Data sets

1) Titanic

2) Iris

3) Digits

According to Data set

- * It can be available for any of the location like big data tools, there will be a remote location to get the data like (Cloud, server)
- * It can be different file format CSV, PDF, XML, JSON, XML

Types of Data

After collecting the data

Generally, ^{tendency of} Data can be Batch data,
Streaming data ↓ Historic Data
↓ Continuous data

1) Structured data → (Table) → ex

2) Unstructured data → (ex: image, video)

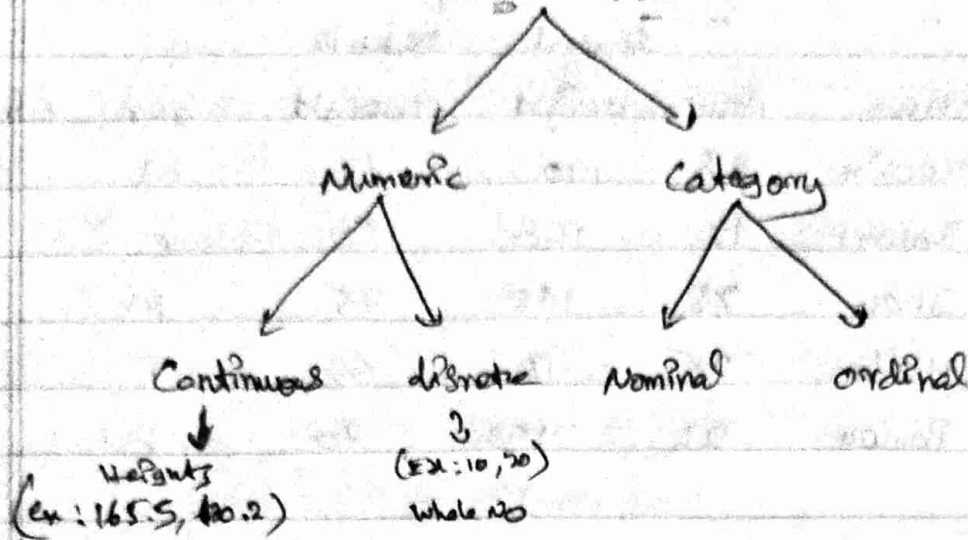
3) Semi Structured data → (xml, JSON)

EDA + FE

Structured Data

feature-1	feature-2	feature-3
weight	weight	Bmi
20	170	22
80	180	24
90	190	26
100	200	30
60	160	27

Structural Data



Category [Nominal, ordinal]

Male
Female } Category

Blood groups } Category

Nominal

* order doesn't Matter that's called
Nominal

Ex: we can have any order of the
data can be fine with the
Data Set.

like male female
female male
male female

↓
This order, it can't be made any
diff in Data Set

ordinal

* order should be important

Ex: 1st → ①
2nd → ②
3rd → ③

Practical Implementation of data

= Student's details

S.no	Name	Age	Height	Weight	Sex	Education
1)	Kevin	28	170	72	M	G
2	Robert	17	140	56	M	HSc
3	John	38	168	75	M	Phd
4	Hellen	25	160	62	F	UG
5	Bailey	24	156	70	P	Phd

First level

* Categorize the features [Columns]

Name	Age	Height	Weight	Sex	Education
↓	↓	↓	↓	↓	↓
Categorize	Numerical	Numerical	Numerical	Cate	Cate
↓	↓	↓	↓	↓	↓
(Nominal)	(Continuous)	(Continuous)	(Cont)	(Nominal)	(Nominal Ordinal)

= Stat Analysis

UNIVARIATE } → single column

BIVARIATE } → two column

MULTI } → more than two column

HSc - ①
UG - ②
PG - ③
Phd - ④

1) I want to check Age's of the students

Univariate Analysis

2) I want to check Height with respective age of the person

Bivariate Analysis

3) I want to check sex with Age & Height

Multivariate Analysis

Independent/Dependent

1) If we want only -

- Age, name, height means it is Independent

2) If we want weight of those people means we are depending on people's age, height so it is dependent.

One Pipeline of ML

- | | | |
|--------------------------------------|---|-------------------|
| 1) Data Ingestion | ↓ | 2) Data |
| 3) EDA | ↓ | EDA |
| 4) Pre-Processing (FFRS) | ↓ | 1) missing values |
| 5) Model Training | ↓ | 2) outliers |
| 6) Evaluating the Model (validation) | ↓ | 3) Scaling |

EDA is

Change

Exploratory Data Analysis is an approach to analyse the data using visual techniques. It is used to discover trends, patterns, or to check assumptions with the help of Statistical Summary and graphical representations.

EDA

Name	Age	Education	Salary	exp
A	23	IT	20k	2
AB	25	BE	25k	2
AC	22	BSc	18k	1
AD	20	BEP	15k	1

EDA Analysis

- 1) Profile of the data
- 2) Statistical analysis
- 3) Graph based analysis

Profile of data

examples

1) # Rows

2) # Column

3) # missing values

4) # numeric values

5) # Column values

6) Data types

7) # duplicates

8) How much data occupied

↓ Statistical Analysis [univariate, bi-variate, multi]

1) Variance

2) cov

3) SD

4) correlation

5) Chi-square Test (or) T-test, (or) Z-test or

6) Mean, median, mode

Anova test

2) Graph based analysis

1) Box plot

2) Histogram

- 3) Scatter plot
- 4) Pie chart
- 5) KDE (Kernel Density Estimation)
- 6) Count Plot
- 7) Heatmap

* Based on FDA we can do a Processing of the Data.

Pre-processing (FF)

- 1) Handle the missing values
- 2) Outlier handle
- 3) Scaling of data
- 4) Transformation (log, Boxcox, Square, cube)
- 5) Encoding
- 6) Imbalance Data
- 7) Feature Selection
- 8) Feature Transformation (Dimension Reduction (PCA, tSNE))

FDA \rightarrow FF

missing value (null values) \rightarrow handle the missing value

outliers deletion \rightarrow handle the outliers

categorical (M, F) \rightarrow encoding (0, 1)

skewed (high range) \rightarrow Scaling the data in certain value

Count of feature \rightarrow handle Imbalance & balance data

If we have sub set of feature \rightarrow Select the particular Some of feature is called feature selection

Dimension reduction (PCA, tSNE)

two feature $x_1, x_2 \rightarrow$ one feature (X)

* We can reduce the Dimensions.

EDA. Automated tool in Python

1) Pandas Profiler

2) mlto

3) Unine

4) Swastix

5) AutoViz

6) Datale