

Types of data : Batch data, streaming data.

minibatch.
data.

(little move.
- frequent)

Historic.
data.

(Periodic).

Continuous data

LIVE.

(Machine Learning)

→ ML

1. Structured data.

→ table

2. Unstructured data. → videos, images,
voice, text.

3. Semi-structured data.

DL.
(Deep Learning)
xml, json

Structured data

Numeric

Category

Continuous

Discrete

Nominal

Ordinal

(Cont.) (Cont.) (Cont.)

HT.

HT.

BMII.

70

170

22

80

180

24

90

190

26

100

200

30

60

160

21.

feature →

1

2

3.

Student Performance dataset → features.

Name	Age	Height	Sex	Weight	Education
------	-----	--------	-----	--------	-----------

Sunny	25	170	Male	70	GD.
-------	----	-----	------	----	-----

Ankit	30	180	Male	80	PG
-------	----	-----	------	----	----

Priyan.	35	160	Male	60	UG.
---------	----	-----	------	----	-----

Priya	20	150	Female	55	Phd.
-------	----	-----	--------	----	------

Aditi	27	145	Female	58	PG.
-------	----	-----	--------	----	-----

Cat.	Nom	Nom.	Cat	Nom	Cat.
Nominal	Cont.	Cont.	Nominal	Cont.	Ordinal

EDA. → type of data.

UG - 0

PG - 1.

Phd - 2.

Univariate. → One column (feature),
e.g Height.

Bivariate. → Two columns.

e.g Height vs Age.

Multivariate. → >2. columns.

e.g Sex vs Height, Weight.

Independent / Dependent variables.

Age, height, Sex.

Independent variable.

define.

height.

Dependent variable

Dependent on
other variable

Pipeline of DS.

(1.) Data Ingestion.

(2) EDA

(3.) Preprocessing. (-feature engineering)

(4.) Model building.

(5.) Evaluation or validation of model.

-feature-/column.

(1.) Missing value handling

(2.) Outlier.

(3.) Scaling of data.

} involves changing
data.
(feature eng.)

Q : First EDA is required or FE ?

(1)

(2)

EDA
Pre-processing

Eg -
Biryani. (Model)

Row data .

{ 1. Chicken .
2. Rice .
3. Onion .
4. Oil .
5. Spices .

EDA.

1 kg → clean ,

1 kg

1 kg

1 ltr .

250 gm .

→ Preprocess .

Cut

ML in real life

Validation . ←

Model .

Tasting .

Biryani .

Dataset

NAME	AGE	EDUCATION	SALARY	Exp.
Scenny	25	UG	25K	2
Deepak	30	PG	30K	3
Ruchi	40	UG	40K	5
Aman	50	Phd.	50K	10
Shalini	20	UG	35K	1

1. EDA (analysis).

Steps.

i) Profile of the data.

ii) Statistical analysis.

iii) Graph based analysis.

Profile of data.

1. No. of rows.

2. No. of cols.

3. No. of missing.

4. Variable type.

Cat ↲

Num. ↳

5. How many duplicates

6. Data types.

7. RAM used.

Stats based interpretation.

- (i) Variation.
- (ii) Covariance.
- (iii) Standard deviation.
- (iv) Correlation.
- (v) Chi-squared.
- (vi) t-test.
- (vii) Z-test.
- (ix) Anova-test.
- (x) Mean/Median/Mode.
- (xi) Skewness.
- (xii) Kurtosis.

Graph based analysis. (Plotting)

- (i) Box plot. → outlier, distribution.
- (ii) Scatter plot → outlier, linearity.
- (iii) Histogram. → distribution.
- (iv) KDE.
- (v) Count base. → Rows, columns.
- (vi) Heatmap. → Correlation.

* Based on EDA we can do processing
of data.

Preprocessing / Feature Engineering. Steps :

- (1.) Missing data handling .
- (2.) Outlier handling .
- (3.) Scaling of data .
- (4.) Transformation (log, boxcox, square, cube)
- (5.) Encoding of categorical data .
- (6.) Handle imbalance data .
- (7.) Feature selection .
- (8.) Dimension reduction , (PCA, t-SNE)

* Automated tool for EDA in python :

- (1.) Pandas profiling
- (2.) MiHD .
- (3.) knime .