

EDA & feature engineering

Data Science Life cycle:

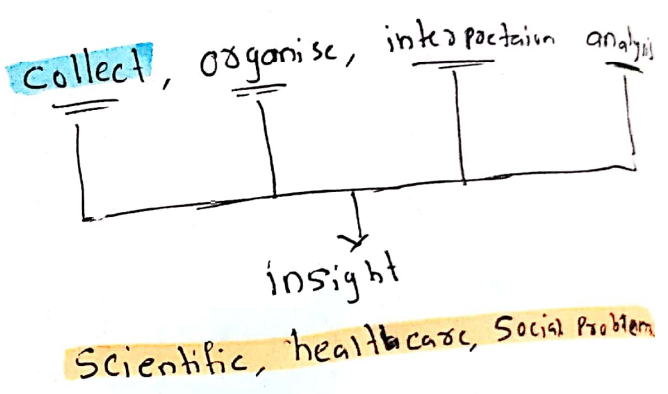
- ① data injection (collection)
- ② EDA (analysis) ← (stats)
- ③ processing (pre) ← (math)
- ④ Model building
- ⑤ Evaluate & validate

machine learning

EDA:

- ↓ Exploratory data analysis
- it often involves the use
- * graphs
- * charts
- * summary statistics
- * and various data visualization

Statistics:



Problem Statement:

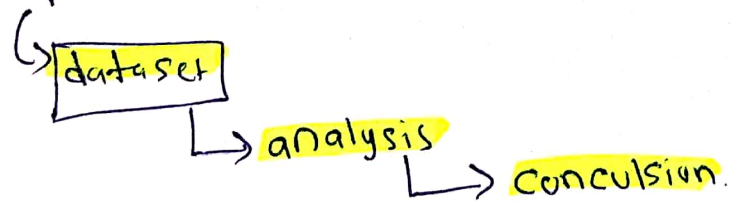
Sales of Product:

→ Sales is going down.

reasons may be

⇒ product, Price, marketing, competitors... etc

We have to prepare a

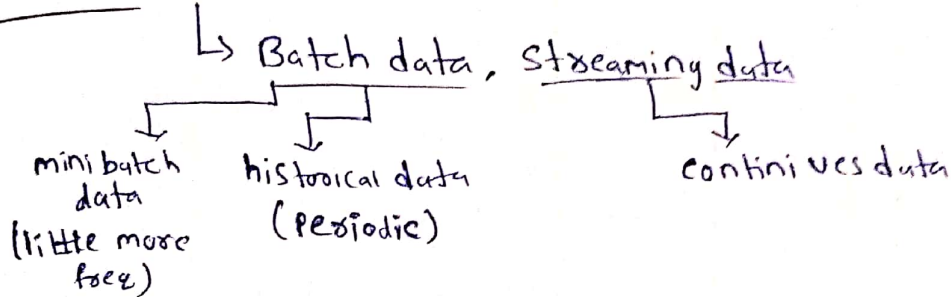


- ① Project manager
 - ② Project Business analyst
 - ③ data scientist
- ↔ domain expert

data can be present in

- hadoop, NOSQL, kafka, spark streaming (HDFS) hadoop distributed file system
- (Bigdata tools), remote location (S3, NOSQL)
- some file formats (csv, tsv, xml, json, excel... etc)
- websites.

types of data:



- ① Structure data :- table format
- ② Unstructure data :- video, image, voice, sound, text
- ③ Semi structure data :- xml, json → deep learning

Structure data:

Weight (kg)	height (cm)	BMI
70	170	22
80	180	24
90	190	26
100	200	30
60	160	21

height is continuous:

like 170, 17.5, 17.65... etc

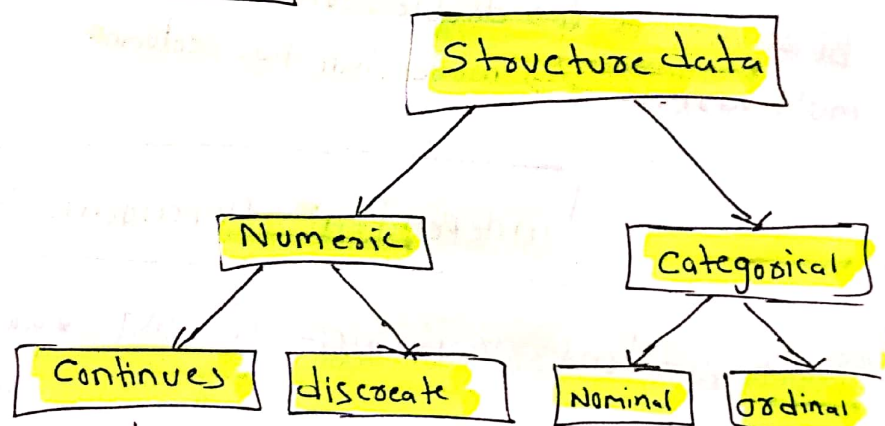
discrete means in a school they are

10, students

20 students

45 students

but we cannot say 25.5, 30.5
so it is standard.



Categorical → male, female, black, white

nominal → order does not matter
red, yellow, black, green

ordinal → meaningful order or rank.
ex: class ranking: 1st class, 2nd class, 3rd class etc.

Ext data sets

assume that

Uni-Variant
bi-Variant
multivariate

Name	Age	Weight	Sex	height	education
Naveen	25	70	male	170	UG
Saikanth	24	60	male	160	PG
Vaasu	35	50	Male	150	UG
Sai	20	70	Male	170	Phd
Baghava	21	76	Male	176	PG
Priya	32	80	Male	180	PG

↑

Categorical
↓
Nominal

↑

Numerical
↓
Continuous

↑

Numerical
↓
Continuous

↑

Categorical
↓
Nominal

↑

Numerical
↓
Continuous

↓

Categorical
↓
(ordinal)

types of data (we have to check)

UG → 0
PG → 1
Phd → 2 } level

Uni-Variate → single column
bi-Variate → double column
multivariate → more than two column.

Independent / Dependent Variable

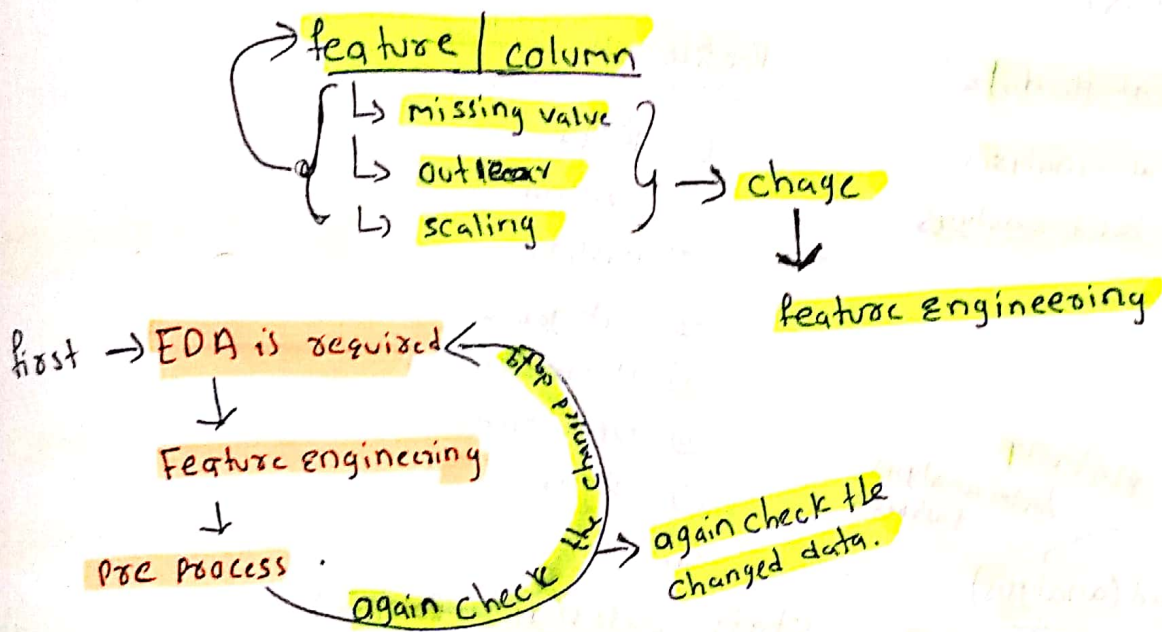
Independent variable - if any variable or value are independent which means not dependent on any value or variable. Called independent ex: Age, Sex

Dependent variable - opposite of independent variable

ex: - weight, height.

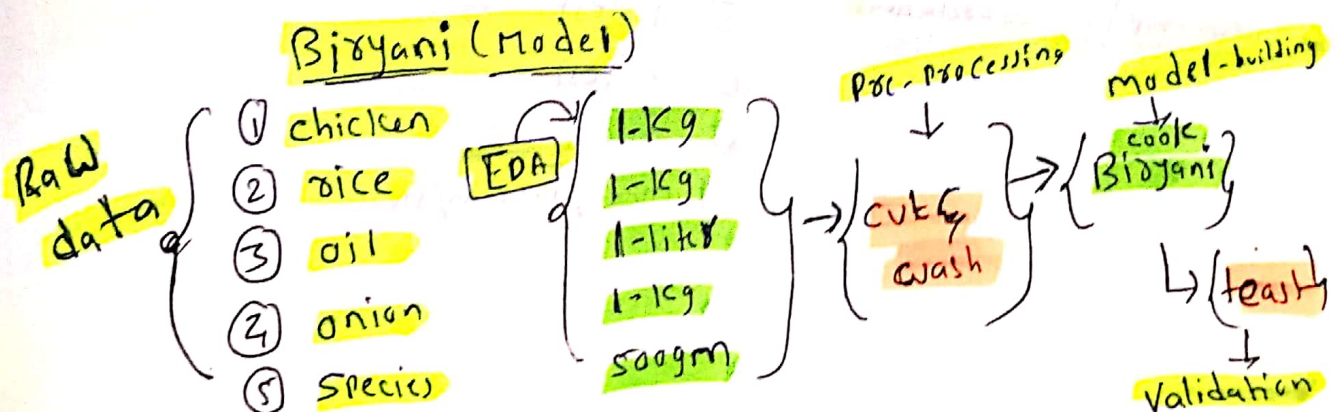
- concept pipeline
- ① Data ingestion
 - ② EDA (Exploratory data analysis)
 - ③ Pre processing \Rightarrow [nothing but feature engineering]
 - ④ model building
 - ⑤ Evaluation or validation of model.

\Rightarrow Data \rightarrow we have to analysis



Practical Example!

SUPPOSE we are doing party and we want to prepare Biryani



⇒ EDA → analysis of data

⇒ Preprocessing (or) feature engineering

we can perform as many times as we wish

⇒
uni-variant
bi-variant
multi-variant

Name	Age	Education	Salary	Emp
Krishna	25	UG	90k	2
Rama	30	UG	80k	3
Hase	40	PG	96k	5
Govinda	50	Phd	87k	10
Jagannadh	20	UG	180k	6

① EDA (analysis)

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

Profile of the data

- ① Row
- ② Column
- ③ missing
- ④ categorical
- ⑤ numerical
- ⑥ duplication
- ⑦ dtype
- ⑧ RAM

Graph based (analysis)

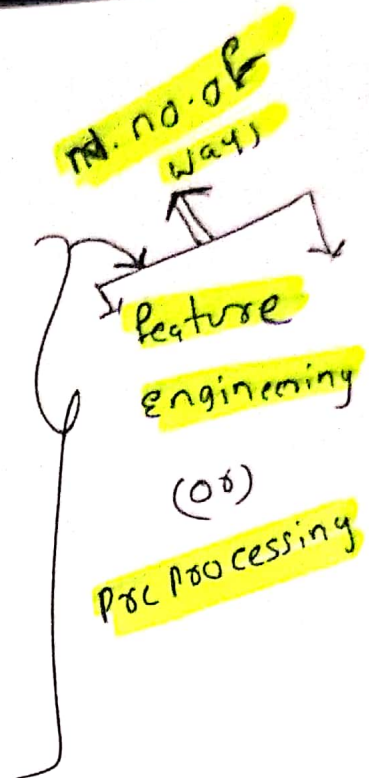
- ① Box Plot → outlier, distribution, statistical conclusion
- ② Scatter Plot → outlier
- ③ histogram → distribution
- ④ kernel density estimation (KDE)
- ⑤ Count bar → Rows, column
- ⑥ heat map → correlation

Stats based (interpretation)

- ① Variance
- ② Co-Variance
- ③ Stand
- ④ Co-correlation
- ⑤ chi square test
- ⑥ t-test
- ⑦ z-test
- ⑧ anova test
- ⑨ mean/median/mode

⇒ Based on a EDA we can do a processing of the data.

- ① missing value handle
- ② outlier handle
- ③ scaling of data
- ④ transformation (log, Box, square, cube)
- ⑤ en-coding
- ⑥ imbalance data
- ⑦ feature selection
- ⑧ dimension reduction (PCA, ESNE)



missing Null value → missing value
EDA → Preprocessing

outlier → handling

categor (man, woman) → encoding

skewed range → Scale (with a certain range)

⑤ → count of feature → handle imbalance
→ feature selection
→ dimension deduction.

Encoding:- Converting the categorical values to numerical value is called en-coding.