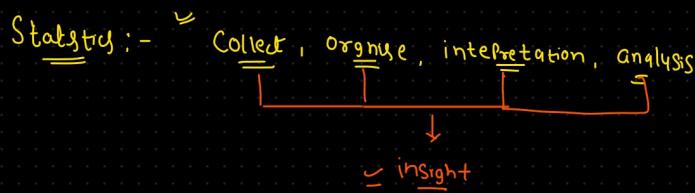
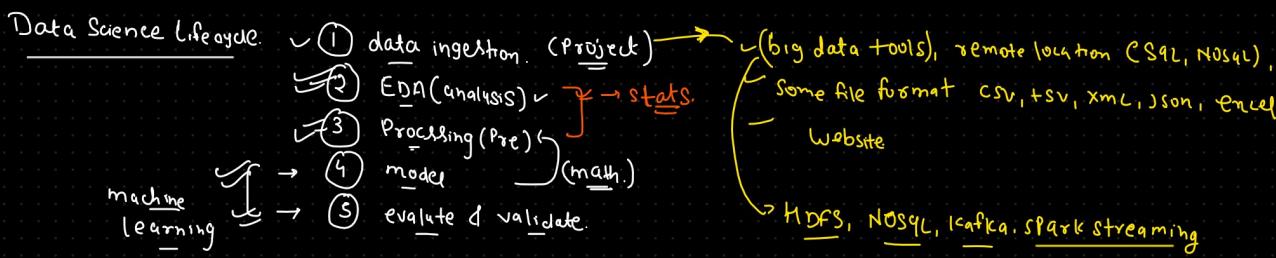


EDA & feature engineering



Scientific, healthcare, social problem

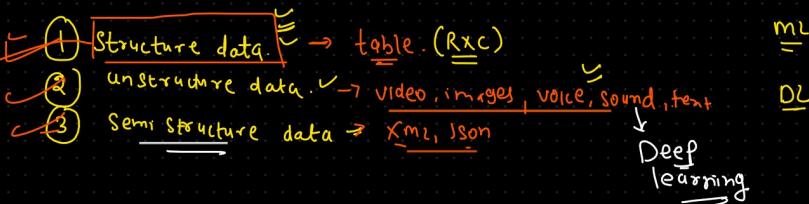
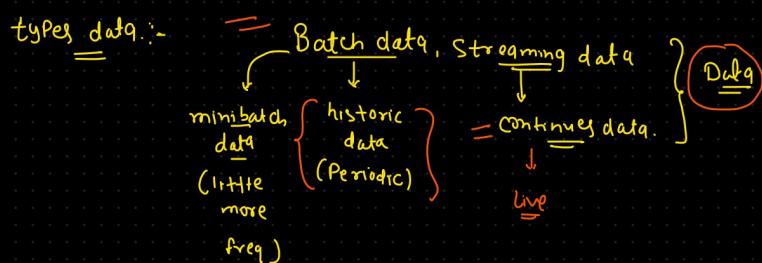
Sales of Product → Sales is going down.

Product, Paying to customer, leadership, marketing, competitor.

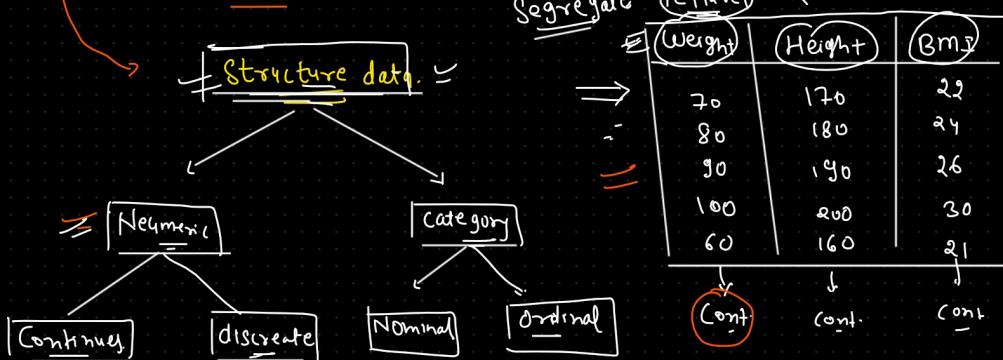
dataset → analysis → conclusion

(1) Project manager
 (2) Business analyst
 (3) Data scientist } → domain expert.

Any domain req. EDA + feature eng.



EDA + FE →



continuum
nature

$$\begin{array}{c} \text{Height: } 160 \underset{\text{=}}{=} 160.5 \underset{\text{=}}{=} 160.55 \\ \quad \quad \quad = \cancel{165} \\ \boxed{[160 \quad 161]} \end{array}$$

$10, 100, 200 \rightarrow$ whole no.

Category → male > cat
female > cat

black > cat
white > cat

- nominal → order does not matter

= male -
= female +

Ordinal → $\begin{cases} 10^{\text{th}} \\ 12^{\text{th}} \end{cases}$
Grad.
→ P.g.
→ phd.

Student Performance

Univariate ↓ bivariate ↓ multivariate ↓ Supervised machine learning

Name	Age	Height	Sex	Weight	Education
Sunny	25	170	male	70	UG
Arijit	30	180	male	80	PG
Priyam	35	160	male	66	UG
Priya	20	180	female	55	Phd
Aditi	27	145	female	58	PG

first-level

✓
25.5
=

nominal
(order
matter)

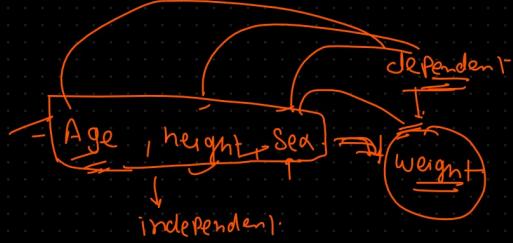
Univariate → single col
bivariate → two w/
multivariate → more than 2

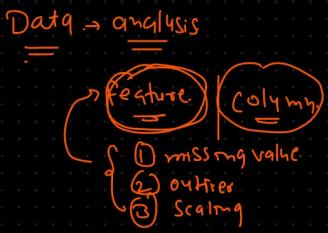
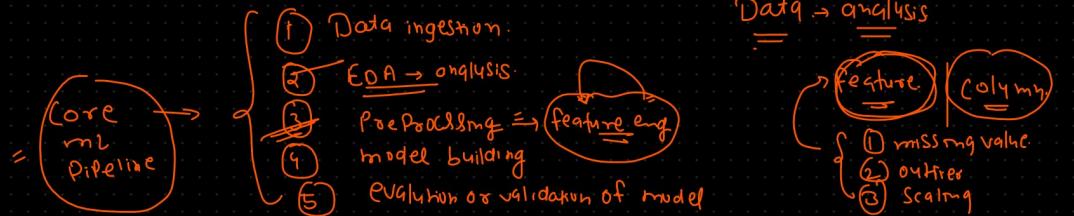
Data.

col/feat

uni, bi, multi

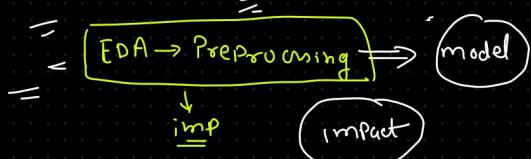
Independent | Dependent



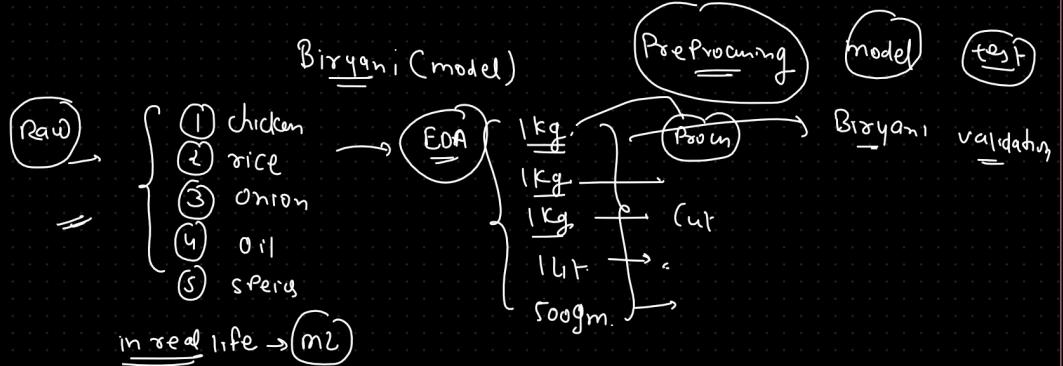


first EDA is req. or FEAT? P? → EDA

change
engineering feature



EDA
Preprocessing



- 1. EDA
- 2. Preprocessing

break till → 4:30

Univariate: Q. EDA
bivariate: Q. Preprocessing Or FE

① EDA ② Preproc.

IID feature
 ↓
 conclusion

- ① EDA (analysis)
- ② Profile of the data
- ③ statical analysis.
- ④ graph based analysis

NAME	AGE	Education	Salary	Emp
Sunny	25	Ug	25K	2
deepak	30	Pg	30K	3
rushi	40	UG	40K	5
Anjan	50	Phd	50K	10
Shalini	20	UG	35K	1

Profile of the data
 { Sunny
 25
 DS
 UG }

- ① Row
- ② Col
- ③ missing
- ④ Cat
- ⑤ num.
- ⑥ duplicate
- ⑦ Dtype
- ⑧ RAM

= Stats based. (Interpretation)

- Uni, bi, multi
- ① Var.
- ② Cov.
- ③ std.
- ④ Correlation.
- ⑤ Chi-square.
- ⑥ t-test
- ⑦ z-test
- ⑧ anova test
- ⑨ mean / median / mode

uni, bi, multi

data analysis
 Plotting
 EDA
 Observation.
 Conclusion

Graph based analysis

- ① Box plot → outlier, distribution, statical, profile
- ② Scatter Plot → outlier, linear
- ③ Pie →
- ④ histogram → distribution
- ⑤ KDE
- ⑥ Countbar R, C (bar)
- ⑦ heatmap → corr

⇒ Based on a EDA. Can be do a Proc. of the Data?

analysis

no. of way

Preprocessing
 of
 data

Feature eng

- ① missing value handle → ✓
- ② outliers handle → -
- ③ scaling of data.
- ④ transformation (log, Boxcox, square, cube)
- ⑤ encoding
- ⑥ imbalance data.
- ⑦ feature Selection
- ⑧ Dim reduction (PCA, tSNE)

$$\left\{ \begin{array}{l} 1) 10 \\ 2) 40 \\ 3) 20 \\ \vdots \\ 4) 10 \\ 5) 20 \end{array} \right\}$$



missing null value → missing value handle

EDA

PP

Outliers → handle

Cat (man, woman) → encoding

Skewed range. → scale (within a certain range)

⑤

Count of feature → handle imbalance

→ feature selection

→ dim red (PCA, t-SNE, LDA)

male
female
0
1

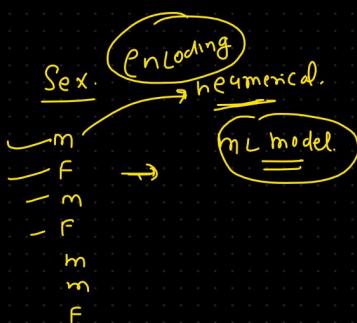
$$y = mx + c$$

$$y = m_{\text{male}}x + c$$

$$0 \times \text{male}$$

$$y = mx_1 + c$$

$$= \text{int}$$



- Practical

EDA

- ① Profile.
- ② State
- ③ Graph.

Data.

Preprocessing

- ① missing → handle
- ② outliers → handle
- ③ scale → handle
- ④ transform → handle
- ⑤ encode → handle
- ⑥ imbalance → handle
- ⑦ drop/duplication → handle
- ⑧ feature selection
- ⑨ dim red

- vast research

Automated in Python

EDA

Pandas Profiling

traits

KNIME

to only 1
with respect to

address 3 automated tools

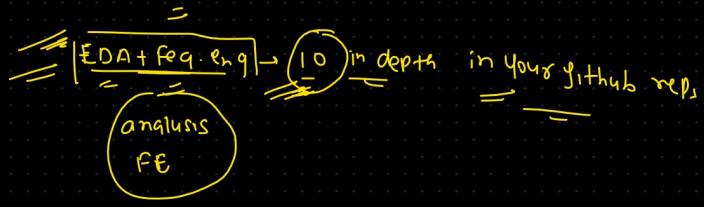
MATLAB

Dataset (10) → 1 → at least 3 automate

max 5 tools

automate.

Pandas Profiling



↓
data
↓
entire eda.

(pp) → again i will write

tomorrow by
12 pm IST

Sunny Squita @ Meuro
h.9i

Krishnaik @ meuron.ai