

EDA & Feature Engineering

Data science life cycle:

- ① data ingestion (Project) \Rightarrow ^{data} (big data tool, remote location, CS21, NO21)
- ② EDA (analysis) \Rightarrow some file format, json, xml, excel, website
- ③ Processing (Pre Processing)
- ④ Model \Rightarrow machine learning
- ⑤ evaluated and validate

Statistics : Collect, organize, interpretation, analysis
Insight

Scientific, healthcare, social problem

Example:

Sales of product \rightarrow sales is going down

* Product, Pricing to customer, leadership, marketing
competitor

dataset \rightarrow analysis \rightarrow conclusion

① project manager

② Business analyst

③ data scientist

} domain Expert

only domain required EDA + Feature Engineering

types data:

Batch data

streaming data

minibatch data
(little more frequency)

historic data
(periodic)

Continuous data

live

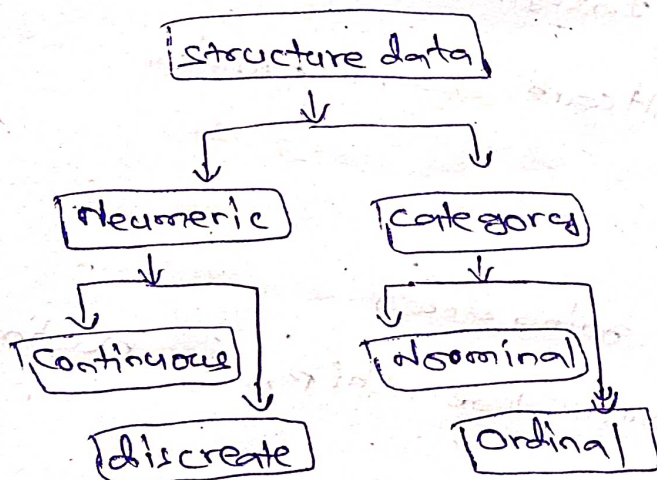
① structure data → table (Row x Column) ML

② unstructure data ⇒ Video, images, audio, text DL

③ Semistrucre data ⇒ XML, JSON Deep learning

* EDA + FE

segregate features



| weight | Height | Bmi |
|--------|--------|-----|
| 70 | 170 | 22 |
| 80 | 150 | 24 |
| 90 | 180 | 25 |
| 75 | 190 | 26 |
| 85 | 155 | 27 |
| 65 | 165 | 28 |
| 60 | 170 | 29 |
| 55 | | |

Continuous: Height 160, 165, 160.55
weight

Continuous Continuous

Discrete: No. of Bank acc 1, 2, 3, 3
children 2, 2, 1, 4

Category: male
female > category

black
white > category

Nominal: order does not matter
male
female

ordinal: order matters

10th
12th
Degree
PG
PhD

Example table

Student performance

| name | Age | Height | Sex | weight | education |
|-------|-----|--------|--------|----------|-----------|
| Ram | 24 | 180 | male | 50 | UG |
| Ajun | 22 | 170 | male | 60 | PG |
| Priya | 23 | 166 | female | 70 | UG |
| Maya | 25 | 150 | female | 80 | Phd |
| Maya | 26 | 160 | female | 65 61 | PG |

↑ ↑ ↑ ↑ ↑ ↑
 Cat Num Num Cate Num Cate
 ↓ ↑ ↓ ↓ ↓ ↓
 Nominal Continuous Contiue Nominal Conti Ordinal

**

Univariate % single column

Bivariate % two column

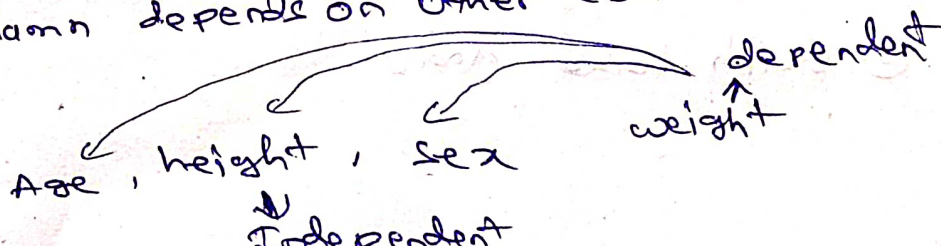
multivariate % ^{than 2} more columns

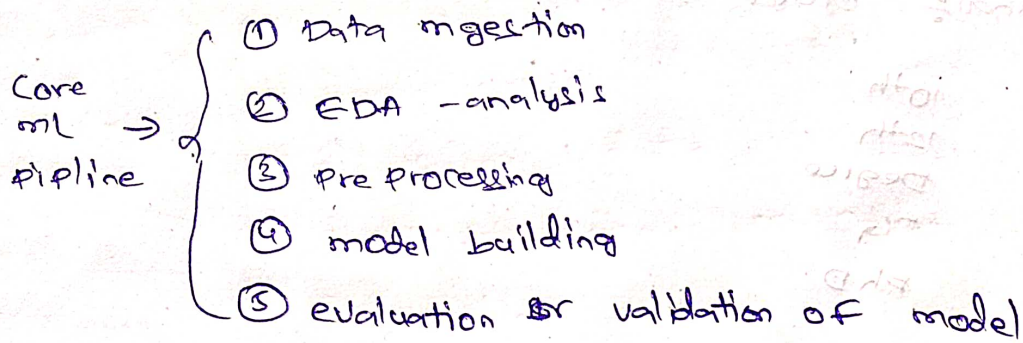
⇒ more than 2 columns

**

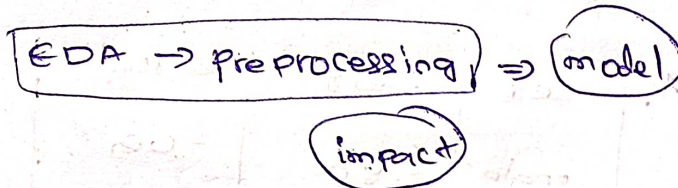
Independent % column doesnot depend on other columns

dependent % column depends on other columns





first EDA is res



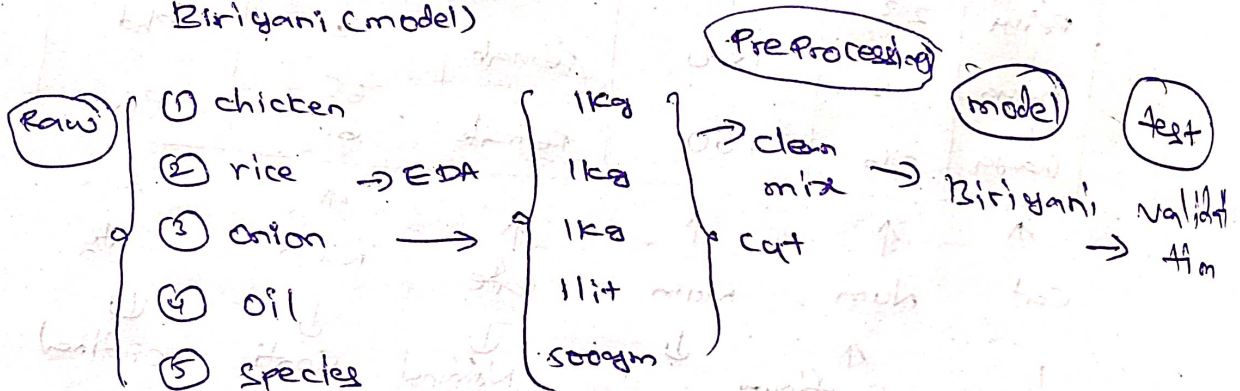
Data analysis

Feature ~~data~~ / Column

- ① missing value
- ② outlier
- ③ scaling

Example % real life

Biryani (model)



① EDA
(analysis)

④ Profile of
the data

② Statistical analysis

③ Graph based analysis

| Name | Age | Education | Salary | exp |
|--------|-----|-----------|--------|-----|
| Ravi | 25 | UG | 25K | 2 |
| Rohit | 30 | PG | 30K | 3 |
| Ravi | 31 | UG | 40K | 5 |
| Kerthi | 40 | Phd | 50K | 10 |
| Shahy | 45 | UG | 25K | 2 |

Profile of the data

Stats based

- ① Row
- ② Col
- ③ missing
- ④ cat
- ⑤ num
- ⑥ duplicate
- ⑦ dtype
- ⑧ RAM

analysis

- ① Var
- ② Cov
- ③ std
- ④ correlation
- ⑤ chisquare
- ⑥ t-test
- ⑦ z-test
- ⑧ anova test
- ⑨ mean / median / mode

univar
bi var
multivar.

Graph Based analysis

- ① Box plot → outliers, distribution, skewed, pos/neg
- ② scatter plot → outlier, linear
- ③ Pie
- ④ histogram → distribution
- ⑤ KDE → R, C (base)
- ⑥ Count bar
- ⑦ heatmap → Cor

using univariate, bivariate, multivariate

dashboard, plotting, data analysis

Preprocessing of data

- ① missing value handle
- ② outliers handle
- ③ scaling of data
- ④ transformation (log, Box-Cox, square, cube)
- ⑤ encoding
- ⑥ imbalance data
- ⑦ feature selection
- ⑧ Dim reduction (PCA, tSNE)

Core ml pipeline

- ① Data collected
- ② EDA (analysis)
- ③ Pre processing or FE
- ④ model building
- ⑤ evaluation matrix & validation

EDA

- ① Profile
- ② stats based analysis
- ③ Graph based analysis

Preprocessing

- ① missing value
- ② Outlier value
- ③ scale
- ④ transformation
- ⑤ encoding
- ⑥ handle imbalanced
- ⑦ feature selection
- ⑧ dimensionality reduction (PCA, tSNE)
- ⑨ duplicate value / duplicate col
- ⑩ split / merge / drop / add

EDA \Rightarrow Preprocessing \Rightarrow model

ways of performing Feature Engineering :

① missing value handle

- ① Random
- ② forward filling / back ward filling
- ③ statistical approach
mean, median, mode
- ④ end of distribution
- ⑤ drop the row
- ⑥ KNN - imputer
- ⑦

③ transformations

box plot

Power transformation

log

square

cube

⑤ Encoding

one hot

label encoding

Binary encoding

target guided encoding

hash encoding

② Outlier

detect

z-score

IQR

boxplot

Scatter plot

Violin plot

handling

drop

median

replace to min-max.

④ scaling

standardization

min-max

unit scaling

⑥ imbalanced

① under sampling

② Over sampling

③ cluster based

over sampling