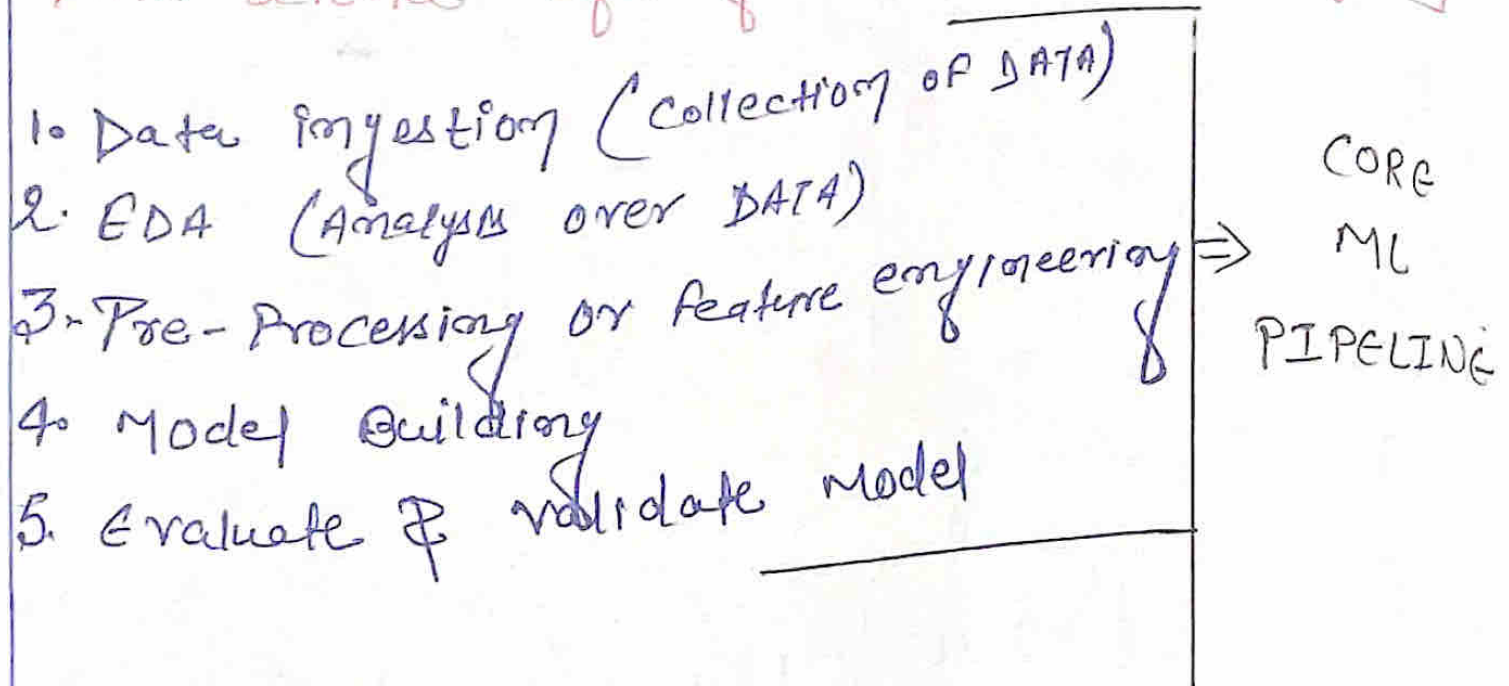


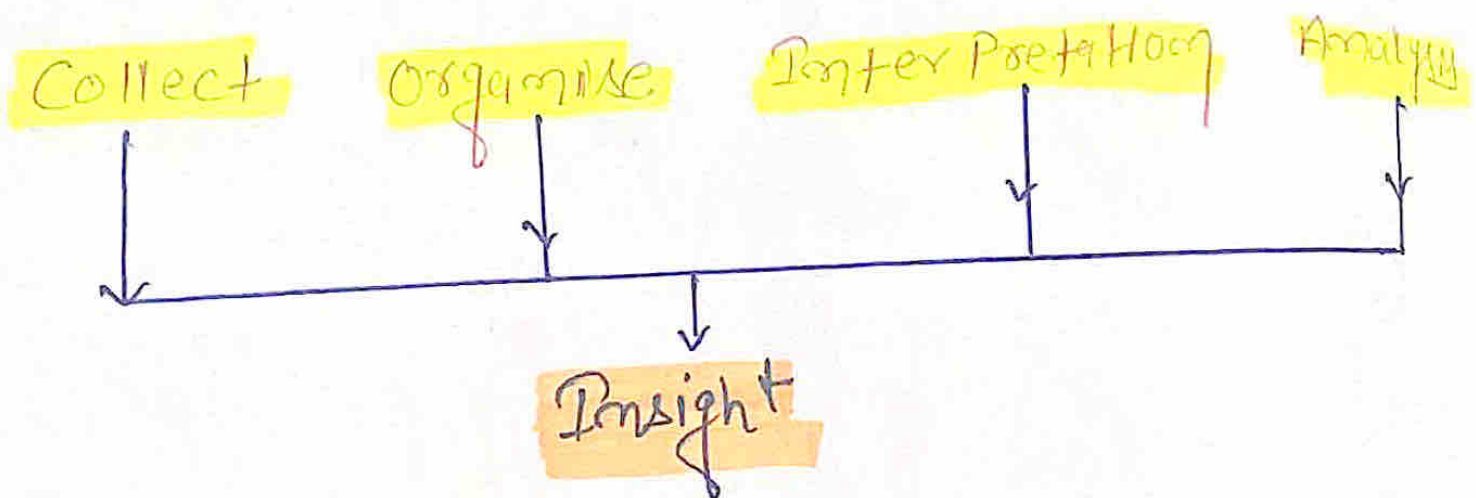
EDA & FEATURE ENGINEERING

DATA Science life cycle: [ML life cycle]



Statistics:

Statistics is the science of collecting, organising & analysing the data



Examples:- Scientific Problem, Healthcare
* every where Statistics use

Problem Statement:-

Sales of Product → Sales is going down

Analysis → Product, Paying Attention to Customer leadership, Marketing, Competitor.

Dataset → Analysis → Conclusion

1. Project Manager
2. Business Analyst → Domain expert
3. Data Scientist → used to get conclusion

Data ingestion: Collecting Data

Examples: Big Data tools, Remote location (SQL, NOSQL), web scrapping

Some file format: (CSV, XML, JSON, xls)

Types of DATA:-

BATCH DATA:- Historical Data, Mini batch Data (Periodic)

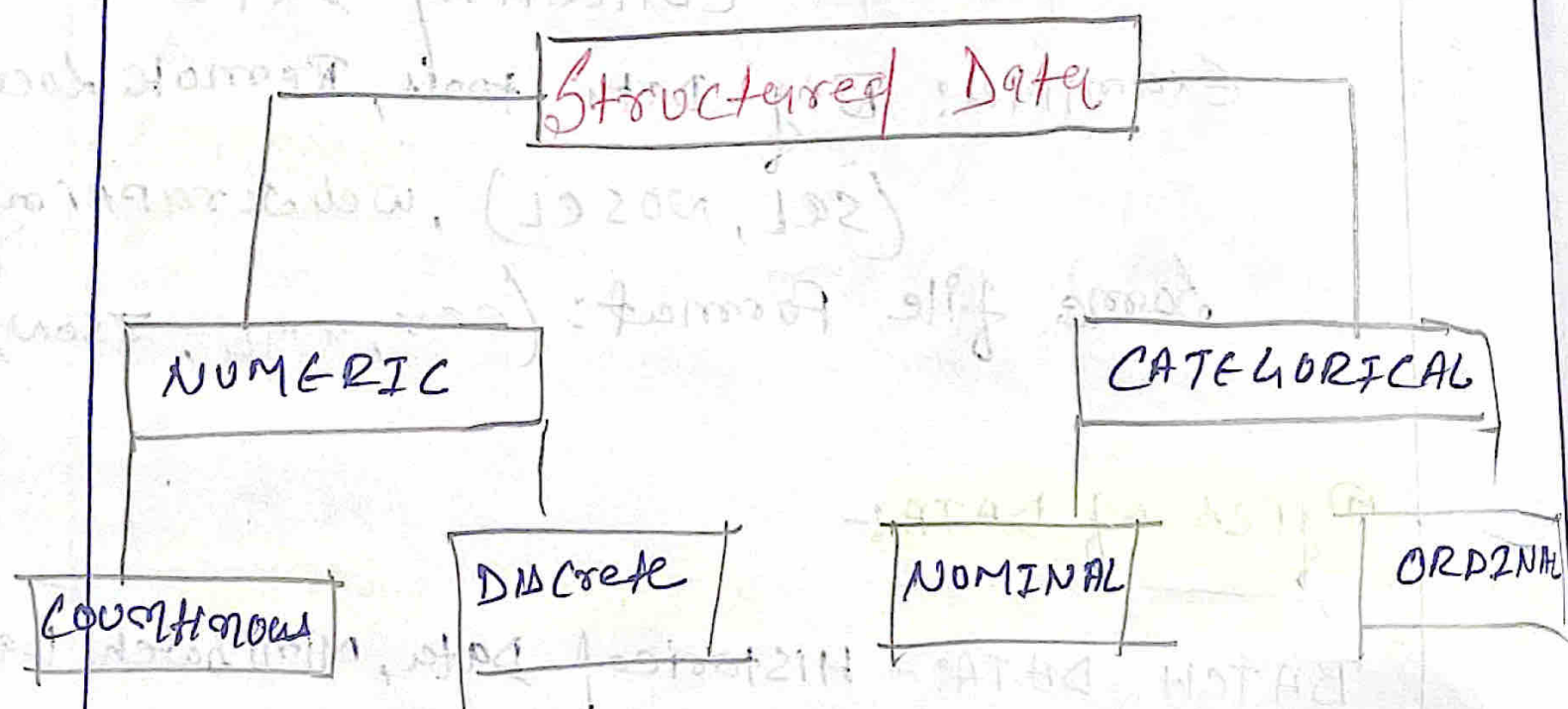
Streaming Data:- Continuous DATA (Live DATA)

1. Structure → Table (Row & Column) → ML
2. Unstructure → video, Images, voice, sound, text → DL
3. Semi-structure → JSON, XML

* Example of Structured Data

Feature 1	Feature 2	Feature 3
weight	Height	BMI
70	170	22
80	180	24
90	190	26
100	200	20
60	160	21

* In Table columns are also called Feature.



NOMINAL :- ORDER Doesnot Matter / MALE/FEMALE

ORDINAL :- ORDER MATTER
Eg:- Degree

10th
12th
UG
PG
PHD

Student Performance (Dataset)

NAME	Age	Height	Sex	Weight	Education
SUNNY	25	170	male	70	UG
ARSHI	30	180	male	80	PG
Priyam	35	160	male	60	UG
Priya	20	150	female	50	PG
Aditi	27	145	female	55	PG
ASHU	35	175	male	85	PG
VISAY	21	140	male	65	PG

Category → NOMINAL
 Age → NUM → COUNT → NOMINAL
 Height → NUM → COUNT → NOMINAL
 Sex → CAT → COUNT → NOMINAL
 Weight → NUM → COUNT → NOMINAL
 Education → CAT → COUNT → ORDINAL

Used to find Data ~~set~~ Types

Univariate → Single column

Bivariate → Two column

Multivariate → More than two column

INDEPENDENT AND DEPENDENT VARIABLE

[Age, Height, Sex] → weight
 ↳ Independent ↳ Dependent

Q. What Required first EDA or FE

⇒ First EDA

Feature engineering

Example:-

RAW
DATA

1. Chicken
2. Rice
3. onion
4. oil
5. Spices

GDA

Preprocessing

1 Kg

1 Kg

1 Kg

1 Litre

500 gm

Biryani

Model

VALIDATION
(TASTE)

* ABOVE is example of ML Model

- ① EDA
- ② Preprocessing
- ③ Model
- ④ Validation

EDA :- Analysis of Data Based on the given Data.

Q. Is Preprocessing & Feature engineering is same?

⇒ Yes, Both are same.

DATASET

NAME	AGE	EDUCATION	SALARY	EXPERIENCE
Sunny	25	UG	25K	2
Deepak	30	PG	30K	3
Rushi	40	UG	40K	5
Amay	50	PHD	50K	10
Shalini	20	UG	35K	1

EDA (Analysis)

- ① Profile of DATA
- ② Statistic Analysis
- ③ Graph Based Analysis

→ Profile of DATA

1. Row
2. Column
3. Missing value
4. Category
5. Numeric
6. Duplicate
7. Dtypes
8. Ram (Size)

⇒ Analysis

Statistical Analysis (Interpretation)

- ① Variance
- ② Co-variance
- ③ Standard deviation
- ④ Co-relation
- ⑤ Chi-Square test
- ⑥ T-test
- ⑦ Z-test
- ⑧ Anova Test
- ⑨ Mean, Median, mode

Univariate
Bivariate
Multi-variate

Graph Based Analysis (Plotting)

EDA

Observation
Matter
for
Conclusion

Also important
Dashboarding

1. Box Plot → Outlier, Distribution, Statistical Profile
2. Scatter Plot → Outlier & Linear correlation
3. Pie Chart →
4. Histogram → Distribution
5. KDE Plot
6. Count BAR ⇒ USED FOR COUNTING (BAR chart)
7. Heat Map ⇒ Correlation

Q. Based on EDA, can we do Pre Processing of Data?

⇒ True

Pre-Processing 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. Dimension Reduction (PCA, tSNE)
9. Duplicate value
10. Split / Merge / Drop / Add (column)

These are steps of feature engineering

1. Missing Null value → Missing value handle (PP)
2. outlier → Handle
3. categorical (man, women) → encoding
4. Skewed Range → Scale (within certain range)
5. count feature }
 - Handle Imbalanced Data
 - Feature Selection
 - Dimension Reduction (PCA, tSNE)

Encoding:- To change categorical Data into numerical Data called Encoding

Some Automated Tools in Python For EDA

1. Pandas Profiling
2. Sweetviz
3. Autoviz
4. D-Tale
5. Mito
6. Kionne
7. DataPrep

DAY 2 OF EDA & F.E

Revision of Day-1:-

Core ML Pipeline

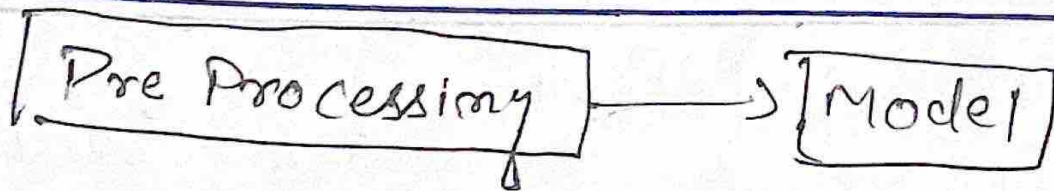
- Data ingestion
- EDA (Analysis)
- Feature engineering or Pre-Processing
- Model Building
- Evaluation matrix or validation

EDA

- ① Profile of the DATA
 1. Statistical Based Analysis
 2. Graph Based Analysis (important)

PreProcessing

1. Missing value handle
2. Outlier handle
3. Scaling of Data
4. Transformation (Log, Cox-Box, Square, Cube)
5. Encoding
6. Handle Imbalanced Data
7. Feature Selection
8. Dimension Reduction
9. Duplicate value / Duplicate Column
10. Split / Merge / Drop / Add (Column)



ways of Performing Feature engineering

① Missing value Handle

- Fill with Random no.
- Forward / Backward Filling
- Statistical Approach (Mean, Median, Mode)
- With the help of end of Distribution, Fill the missing value
- Drop the Row
- ⇒ Impute with KNN (KNN-IMPUTER)
- ML Algorithm for missing value
- Build own ML Model to Predict Missing value.

② Outlier Handling

Detect the Outlier

- Z-Score
- IQR
- Box
- Scatter-Plot
- Violin-Plot

Handle outlier

- Drop
- Fill with Median
- Replace / Impute

⑤ Transformation

- ① Box Cox Transformation
- ② Power Transformation
- ③ log
- ④ Square
- ⑤ cube

4. Scaling of Data

- i) Standardization
- ii) Min Max Scales
- iii) Unit Scaling

⑤ Encoding Various Method

- ① one-hot encoding
- ② Label - encoding
- ③ Binary encoding
- ④ Target Guided Encoding
- ⑤ Hash Encoding

6. Imbalanced Dataset (Treatment various method):-

- ① collect more Data
- ② under Sampling
- ③ over Sampling
- ④ Cluster Based over Sampling

Q. How to Find the Best Model Accuracy
various Method?

→ To increase the Accuracy, we need to
change the Pre-Processing technique
use different method or step from
the above.

→ we need to use each & every Pre-
Processing step & find Best Accuracy.

Q. How do we transform the Data?

⇒ import numpy as np

np.log(df)

sns.distplot(df)