

4) Scaling

- (a) Standardization
- (b) Min-Max
- (c) Unit Scaling

5) Encoding

- (a) One hot
- (b) Label encoding
- (c) Binary encoding
- (d) Target guided encoding
- (e) Hash encoding

6) Imbalanced dataset

- (a) collect more data
- (b) Under sampling
- (c) Over sampling
- (d) ~~over~~ & cluster based

Note: we can change preprocessing technique to improve accuracy.

24-SEP-2022

EDA and Feature Engineering

Data Science lifecycle

- ① Data Integration → ^{Project's discussion}
 - ② EDA (Analysis of data) → ^{Exploratory Data Analysis}
 - ③ Pre-Processing
 - ④ Model/Algorithm Building → ^{Different Algorithms}
 - ⑤ Evaluate & validation of the model
- Use statistics to analyse the data
core ML pipeline

Statistics

→ collect, organize, interpretation, and analysis of data.

we can find some insight by performing these actions.

e.g. Sales of Product → sales is going down.
what is the reason behind this down?

→ Is our product unique?
not paying attention to the customer
Leadership is not good.
Marketing strategy is not good.
Not looking to the competitors.

Using all these we created one dataset
need to analyse to find out the reason.
↓
conclusion.

① Project Manager

② Business Analysts

↓ ③ Data Scientist

Domain Expert

} look into the dataset and find conclusion.

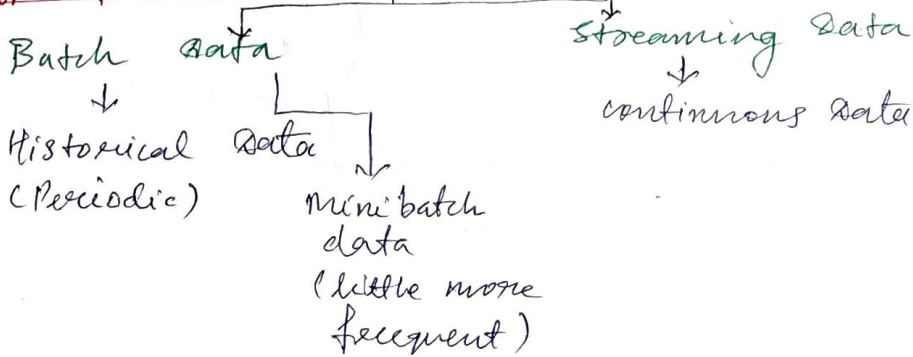
Note: Any domain requires EDA and feature engineering.

Data Integration

- ① from big data tools like Hadoop, HBase, Spark etc
- ② remote location like S3, MongoDB etc
- ③ Some file format like CSV, XML, TSV etc

Types of data

Tendency of Data



- ① Structured data \rightarrow Table (RXC) \rightarrow ML
- ② Unstructured data \rightarrow Videos, images, voice, text \rightarrow DL
- ③ Semi structured data \rightarrow XML, JSON

EDA + FE \rightarrow In perspective of type of data, we can perform.

Structured data

eg

feature 1	feature 2	feature 3
weight	Height	BMI
70	170	22
80	180	24
90	190	26
100	200	30
60	160	21

continuous continuous continuous

Two types of Structured Data

- ① Numerical
 - continuous \rightarrow continuous in nature. eg: Height
 - discrete \rightarrow whole number eg: student in the class
- ② Category
 - Nominal \rightarrow Order doesn't matter eg: Gender (M or F)
 - ordinal \rightarrow Order matters eg: Degree (10th, 12th, Graduation, Master, PhD)

Implementation of the Data

Student Performance

Name	Age	Height	Sex	Weight	Education
Sunny	25	170	Male	70	GDVU
Arijit	30	180	Male	80	Ph
Priyam	35	160	Male	60	U
Priya	20	150	Female	55	PhD
Aditi	27	145	Female	58	Ph

↓
 Categorical
 Nominal
 Ordinal

↓
 Numerical
 Discrete
 Continuous

↓
 Numerical
 Continuous

↓
 Categorical
 Nominal

↓
 Numerical
 Continuous

↓
 Categorical
 Ordinal

{ U → 0
 Ph → 1
 PhD → 2 }

① First-level

① Categorize the feature and check type of data

② Second-level

Univariate } → one column → check Height column only.
 Bivariate } → Two column → check Height w.r.t Age column.
 Multivariate } → More than two column → check sex w.r.t Height and Age column.

Independent and Dependent Variable

Age, Height, Sex → Based on it defined weight
 Independent feature
 Dependent feature

→ Independent means independent analysis of the feature.
 → Dependent means one feature is dependent on one or more feature.

Pre-Processing Feature Engineering

① Missing values
 ② Outlier detection
 ③ Scaling

} Changes in the column → Transformation
 next feature

→ Engineering w.r.t feature.

Ques. What EDA is required or F.B or P-P?

Ans. EDA → Pre-Processing → Impact model

EDA and F.E

e.g. Real life example -
Party
Cook Biryani (model)

- ① chicken
 - ② Rice
 - ③ Onion
 - ④ Oil
 - ⑤ Spices
- } Raw Data
to cook chicken

1 kg
1 kg
1 kg
1 kg
500 gm

} → Directly we can't cook Biryani.
we have to clean oil, cut onion etc as a preprocessing steps.

↓
Cook Biryani (model)

↓
Taste Biryani (evaluation and validation)

EDA → analysis of data

Preprocessing or F.E → cleaning the data.
Preprocessing and Feature Engineering are same.

→ Ist we will perform EDA, then F.E after that we can again perform EDA.

e.g.

<u>Name</u>	<u>Age</u>	<u>Education</u>	<u>Salary</u>	<u>Experience</u>
Sammy	25	UG	25K	2
Deepak	30	PG	30K	3
Rishi	40	UG	40K	5
Aman	50	PhD	50K	10
Shalini	20	UG	35K	1

EDA → Analysis of data

- ① Create profile of the data
- ② Statistical based analysis
- ③ Graph based analysis

Profile of the Data

- ① Number of Rows
- ② Number of columns
- ③ How many missing values
- ④ How many numerical values
- ⑤ How many categorical values
- ⑥ What is the type of the data
- ⑦ Duplicate value
- ⑧ How much RAM is consumed

Statistic based interpretation

→ univariate, bivariate, multivariate

- ① Variance
- ② Covariance
- ③ Standard deviation
- ④ Correlation
- ⑤ Chi square
- ⑥ t-test
- ⑦ z-test
- ⑧ F-test
- ⑨ mean/median/mode
- ⑩ skewness
- ⑪ Kurtosis

Graph based analysis

- ① Box Plot → can check outlier, distribution.
- ② Scatter Plot → can check linearity, outlier
- ③ Py Plot
- ④ Histogram → distribution of the data
- ⑤ KDE plot
- ⑥ Q-Q plot
- ⑦ Heatmap → correlation b/w variables
- ⑧ Countbar → Count Row, column for each variable

Ques Based on a BDA can we do a preprocessing of the data?

- Ans Yes, we can do preprocessing ~~stop for BDA~~ ^{and}
- ① Missing value can be handle
 - ② Outlier can be handle
 - ③ Scaling can be done
 - ④ Transformation is possible

- ⑤ Encoding is possible
- ⑥ can handle imbalance data
- ⑦ feature selection is possible
- ⑧ Dimension Reduction (PCA, tSNE)

Missing value detection \rightarrow Missing value handle
EDA Pre-Processing

outlier detection \rightarrow outlier handling
EDA Pre-Processing

cat (man, woman) \rightarrow encoding
EDA Pre-Processing

skewed range \rightarrow Scale (within certain range)
EDA Pre-Processing

weight feature \rightarrow handle imbalance data
EDA Pre-Processing

\rightarrow Feature selection

\rightarrow Dimension Reduction (PCA, tSNE)
Pre-Processing

① 10
② 10
③ 20
④ 10
⑤ 20
} Imbalance dataset

Encoding

$$y = mx + c$$

male, female

$$y = 0 \times \text{male} + c$$

not making sense
need to encode male and female to 0,1 so
that machine can understand.

Note: In EDA, we do analysis, and in Pre-processing or Feature Engineering we do change the value.

Summary

BDA

- ① Profile
- ② Statistical Analysis
- ③ Graphical Analysis

Preprocessing (directly effect model)

- ① Missing values
- ② Outlier
- ③ Scaling
- ④ Transform
- ⑤ Encode
- ⑥ Imbalance
- ⑦ Drop/Amplify
- ⑧ Feature Selection
- ⑨ Dimension Reduction (PCA, LDA, tSNE)
- ⑩ split/merge/drop/add

Automated tool in Python

- ① Randa's Profiting
- ② mito
- ③ Krime

[Note: Before giving data to the model, we do EDA and Pre Processing.]

Pre-Processing and Feature Engineering ways

① Missing value Handle

- ① Random
- ② forward filling / backward filling
- ③ statistical approach,
 - mean
 - median
 - mode
- ④ end of the distribution
- ⑤ drop
- ⑥ KNN-Imputer or diff imputation technique
- ⑦ ML algorithm which supports missing values
- ⑧ we can create our own ML model and can predict missing values

② Outliers handle

(i) Detect outliers

- ① Z-Score
- ② IQR
- ③ Box Plot
- ④ Scatter Plot
- ⑤ Violin Plot

(ii) Handling of outliers

- ① Drop
- ② Replace ^{Fill} with median
- ③ Replace ~~with~~ trimming

[Note: Does outliers affect mean?
→ Yes it affects the mean.]

③ Transformation

- ① Box-Cox
- ② Power Transformation
- ③ log-Transformation
- ④ Square
- ⑤ cube
- ⑥ Yeo Johnson