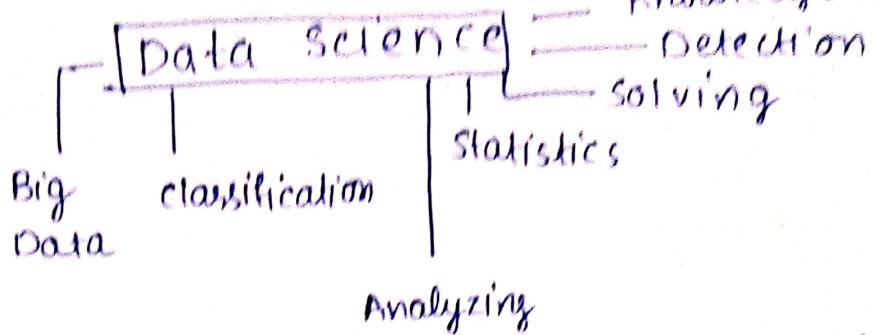


UNIT - 1

Q. what is Data Science?

Data = collection of information
 Row = horizontal representation of data
 Column = vertical representation of data
 Knowledge = Data



- understand patterns, trends, correlations with in data
- make predictions. (scientific words)
- Recommendation based on data. Yhi we honge
- make informed decisions { Data science will able to make }
 - better decision
 - what will happen next
 - pattern finding
- Solve Real world Problems

Data science multidisciplinary field hoti hai jo Data pr scientific methods use Karke process Karti hai algorithm & system ki help se Extract Karti hai, knowledge & insights structured/unstructured data see

Applications of Data Science

- (1) Healthcare - (a) pattern dekh ke pta lagate hai ki Future mein Konsi disease fastly spread hogi.
 (b) optimization ke sath treat Karti hai patient ko like past data set ki help se deka jata hai yeh symptoms konsi disease ke the ya Konsi disease us person ko past mein hui thi.
- (2) Finance - (a) stock market mein Kya Bachhega kya Ghatega
 (b) fraud detection Karta hai.

(3) E-commerce - (a) past purchasing/Interest Data dekh ke Recommend Karta hai

(b) product demand batata in future increase hoga ya decrease.

(4) Marketing - (a) social media se customer ka interest dikھata hai.

(5) Education - (a) Decision making sikھata hai.

(6) Human Resources - (a) forecast Karta hai Business mein aage Kya hoga.

Data Scientist

SKILLS NEEDED

- solid background in mathematics, statistics, programming & domain specific knowledge.
↳ Python, R (technical)
- Data related tasks like cleaning, modelling machine learning algo Kaise apply karte hai. (technical)
- communication. (non technical)
- Business understanding (non technical)
- critical thinking approach. (non technical)
- Machine learning & AI (technical)
- technical (team work/collaboration)

Roles & Responsibilities of Data scientist -

Management -

Supporting the planning & Executions

Analytics

Create models for analytics problem.

Strategy

Design strategy by understanding customer's trends & behavior.

Collaboration -

collaboration with people like stakeholders.

Provide business solution using data Related techniques.

Create AI tools to utilize data.

Recommend ways to improve Data Reliability Efficiency Quality.

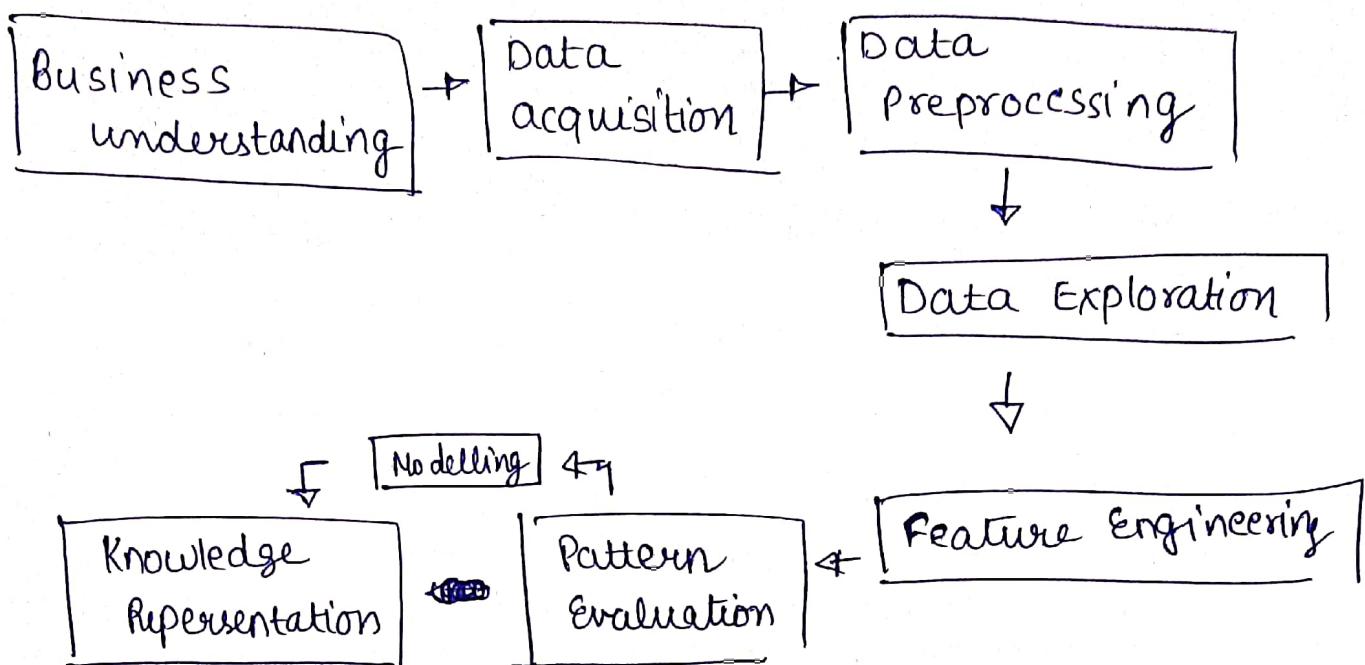
Automate work through the use of predictive analysis

Explore & Examine data find hidden patterns

Why Python used in Data Science -

- (1) interpreted language, code is executed immediately.
- (2) flexible programming language.
- (3) Easy & Free to use.
- (4) have large community support.
- (5) Extensive supports of libraries.
[Numpy, Scipy, Pandas, Matplotlib & Seaborn]
- (6) Built in mathematical functions.
- (7) versatility - allow data scientist to perform a wide range of tasks.
- (8) Data handling - cleaning, manipulation & handling missing values main help Karti hai.
- (9) Integration - Semlesly integrate Karti hai other data science tools & technologies ke sath or work Karti hai.
- (10) Indentation.
- (11) Portable
- (12) Straight forward syntax.

Life cycle of Data Science



Business understanding - problem ko define Karo, (statement)

Ki problem mein Karna kya hai, Kis cheez ka solution nikalna hai and kiske lie kya Data set needed hai.

Data acquisition - Data ko collect Karo jo problem se Related hai, or yeh Data boht sare sources se atta hai like ko .pdf, .image jpeg/ jpg, .esv etc. So, ism sb ko ek hi format mein convert Karna.

Data preprocessing - Data Preprocessing Ek process hai, jisme Apne kaam se Related Data ko leta hai or Baaki sb ko Remove kar dete hai.

Data Exploration -

is phase mein Data ko explore karte hai, ki the rows, kitne columns, kya values, Mean, median, variance & all kitna hai.

Feature Engineering -

Python tools ki help se missing values and errors ko treat karte hai jisse jo predictions ho vo thik aye.

Pattern Evaluation -

pattern Evaluation mein hum apne result ko evaluate karte hai, ki jo predictions ki hai vo shi hai ki nhi, nhi hota to vapis se data ko process karte hai or Ajar result shi hota hai to next step follow karte hai.

Modelling -

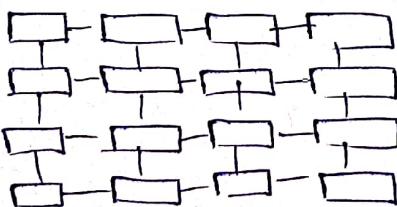
Apne predictions ko model pe apply karte hai.

Deployment -

Model present Karenge shareholders ko or feedback lenge. jo Kami hongi vo shi Karenge or market mein model ko hunch launch Karenge.

Structured Data

- Data follows well defined, and organized structure, yeh store hota hai databases mein tables & rows ki form mein.
- organized hota hai Rows & columns ki form mein.
- tabular formats of Relational databases, ki form mein hota hai.
- less flexible
- Efficient for analytical processing.
- Business transactions, financial records.
- Data nature - Quantitative
- storage - Requires less storage
- analysis - Easy
- Graphical Representation.



Unstructured Data

Data lack karta hai predefined model or structure yeh kai baar images, videos ki form mein hota hai.

Koi predefined structure nhi hota.

stored in various formats

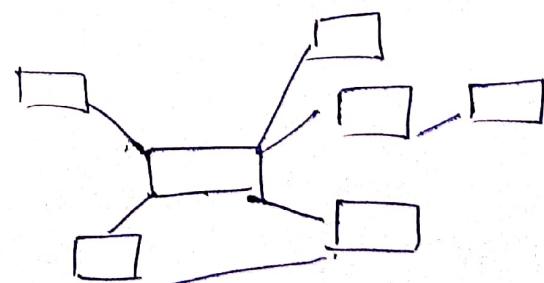
More flexible

processing may require advanced technique like natural language processing.

Speech & image Recognition

Qualitative

Storage - Requires More storage different.



techniques to process unstructured Data -

- ① Tokenization - Breaking down → text ko break karne small words or phrases mein.
- ② Part-of-speech tagging - Assigning grammatical categories (noun, verb, adjective) noun, verb, adjective add kardo sentence mein.
- ③ Sentiment analysis - Determine karna anyisi's karna tone ka, tone kaisi hai (Positive, Negative ya Neutral).
- ④ topic modelling - identify karna themes or topics ka jo document ke hote hain q'like Related kis cheez se hai}
- ⑤ Text classification - text ko categories mein divide karna For EXP - Email $\begin{cases} \text{spam} \\ \text{not spam} \end{cases}$

challenges in processing unstructured data -

1. Lack of structure - organized nahi hota islie analyze karne main problem hoti hai.
2. volume & variety - Kai baar data kaifi Badli Quantity mein hota hai, islie or kaifi different formats mein bhi hota hai islie problems aati hai analyze karne mein.
3. Data noise - irrelevant information/mistakes (in the form of typos, abbreviations or irrelevant details)

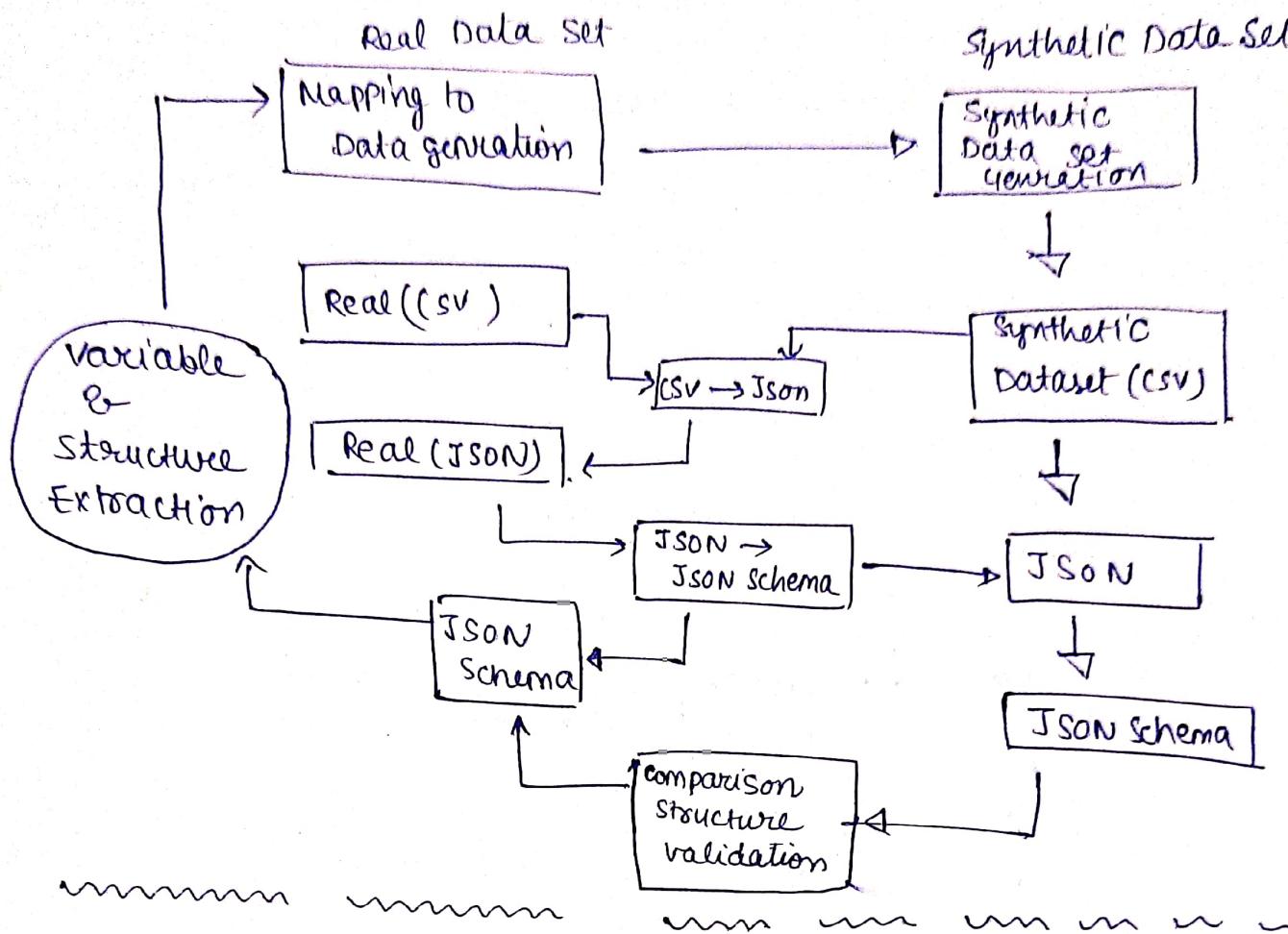
Data integration - storage methods alag hote hain
structured or unstructured ke.

Scalability - scalable solutions chahiye hote hain
data ko process karne ke lie.

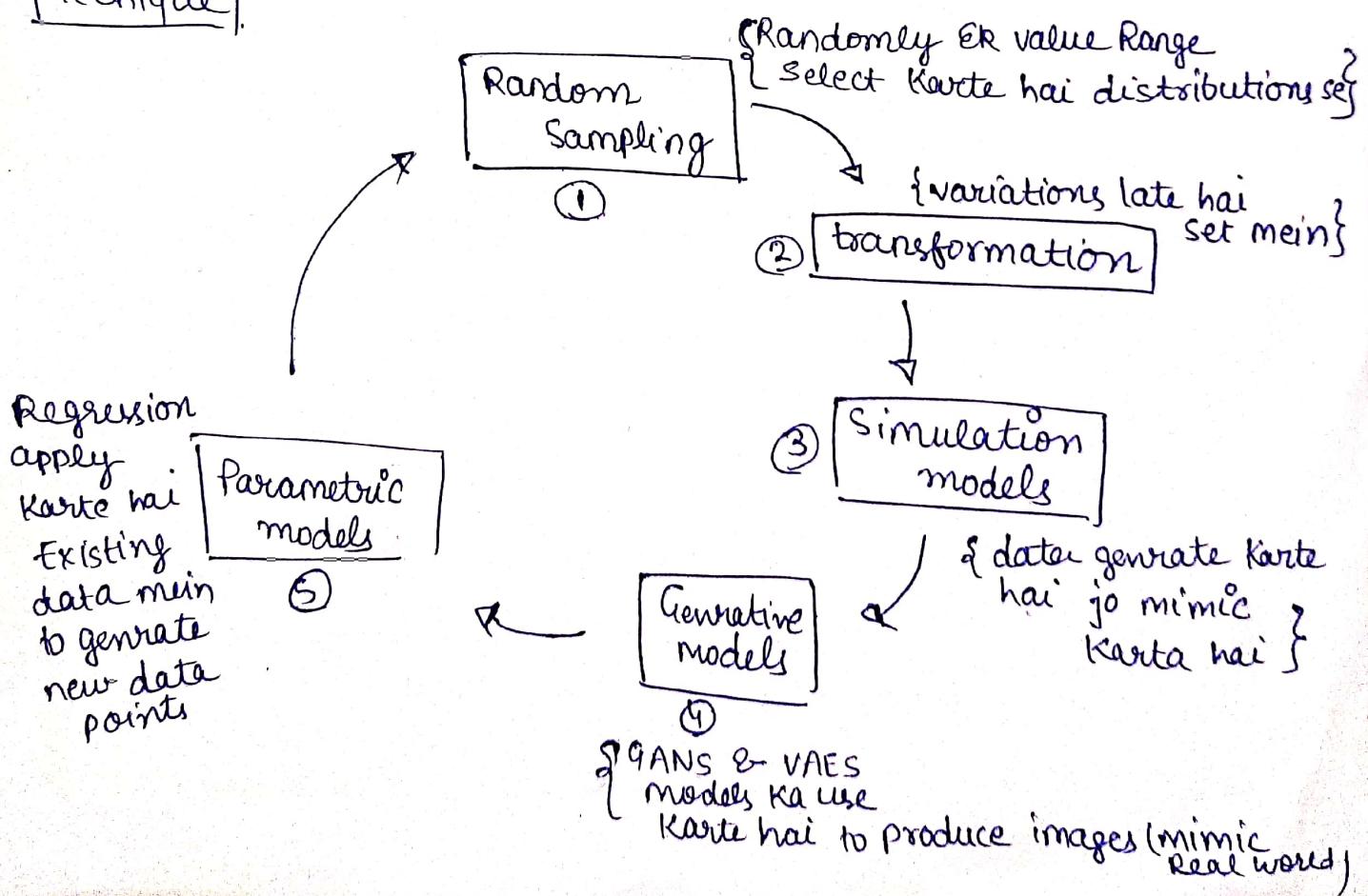
Privacy concerns - data sensitive information carry
karta hai, so, privacy ka bhi dhyān
rakhna padta hai.

continuous Evaluation - unstructured data kaafi quantity
mein generate hota hai, so, tools chahiye hote
hai large unstructured ko process karne ke lie.

CSV → JSON Synthetic Data Set



Technique



Privacy in synthetic Data generation

→ Synthetic Data generation involve Karta hai creating artificial Data jo mimic Karta hai Real Data set ko.

① Ethical considerations - Respecting individual privacy is an Ethical Responsibility. It Should not compromise the privacy or confidentiality of individuals represented in original Data set.

② Building trust - maintaining privacy in synthetic Data generation Fosterst among stakeholders, Encourage Karta hai unko participate Karne ke lie data share karne ke lie bina presham (concern) hue.

③ Mitigating Risks - privacy violation hoga, Synthetic Data mein sensitive information provide Karana jiske lie jail thi ho skti hai.

④ Protection Against threats - koi thi kisi thi information ka misuse na kare islie synthetic Data avoid Karta actual sensitive information ko share Karne

⑤ Preserving long term Data utility - original Data Set outdated ho jate hai Regulation, Technology or Data Quality issues ki wajh se . jiski privacy controlled synthetic Data kaafi time the sustainable Rہتے hai.

It can be done by - {mask like jaise face mask se face chupate hai vo wale}

① Anonymization - mask Kar dete hai personal information ko

② Perturbation - introduction of small/controlled changes or disturbances in data set Jisne Data set ki Prediction affect nhi hoti.

③ Data Generalization - Aggregate Karna / generalize Karna Data ko hide Karne ke lie fine grained details ko.

④ Differential privacy - noise add Karna records mein re-identification ko prevent Karne ke lie.

UNIT - 2

Python Basics

* constants - fixed values like numbers, strings

Print(82)
 └─constant
Print('I love python')

* Reserved words - Predefined words jo as a variable use
nhi kie ja skte, Ex - true, False, finally, for,
try, except, import etc.

* variables - define karta hai information ko memory mein
jo as a fetch Baad me use ki ja skti hai

a = 23
b = 'Python is great !!'

* Assignment - assigning value to variable

a = a + 5

* Numerical & Numerical expression -

- + → Addition
- → Subtraction
- *
- / → Division
- ** → Exponentiation
- % → Remainder

for solving precedence

- Parenthesis
- Exponentiation
- Remainder, multiply, divide
- Addition, subtraction
- Left to Right

* Boolean & comparison - have two possible values (T/F)

comparison operator

Less than $\textcircled{a} <$

Less than Equal to $\textcircled{a} \leq$

Greater than $\textcircled{a} >$

Greater than or Equal to $\textcircled{a} \geq$

Equal to $\textcircled{a} ==$

Not Equal to $\textcircled{a} !=$

* Strings - created by using single or double quotes
string - 1 = "I am valid!".

concatenate using + operation.

length can be find by len(string - 1)

.title make the Every^{1st} word capital.

* Type . - Integer

- float

- string

- Boolean

* functions - can write code in blocks, It makes code more Readable & Easy to Reuse.

amount = hours * Pay - per hour

Hum isko multiple locations per use Karna chalte hain
to isko copy paste kar skte hain

but ek din amount mein 500 add karna tha.
to formula

amount = hours * Pay - per - hour + 500

ab code mein yeh line dhund ke edit karne padgi
islie yeh sb baar na karna pade,
function introduce kiya.

→ used to define a function

def calculate - pay (hours, Pay - per - hour):

 amount = hours * Pay - per - hour

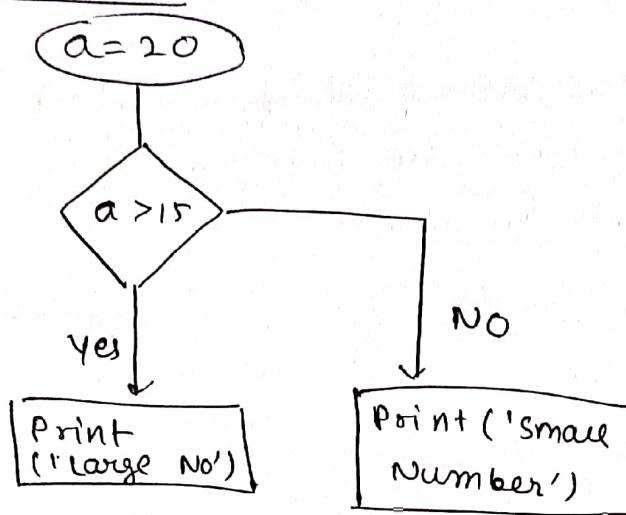
 return amount → to return a value

values passed
to function

now jab bhi hume is formula ki zaroorat hogi to
call karenge function ko.

calculate - pay (20, 120)

Conditional statement -



Syntax

if <condition>: conditional Expression

 do - something
 ↓

True hai to aage badho execute karo
varna skip karo

14 tab spaces

if Else

```

Price = 200
if Price < 150:
    Print("Buy !!")
else:
    Print("sell !!")
  
```

if Elif - Else

```

Price = 200
if price < 150:
    Print("Buy")
elif price >= 150 and price < 250:
    Print("Hold")
else:
    Print("sell").
  
```

and → dono statement ko check karta hai "dono true hone chahiye"

OR → I bhi true hoti hai to statement execute hoti hai.

not → If koi statement true hogi to false bolenge or false hogi to true.

Data structures & loops

List - Data type.

Values same type ki honi chahiye

[] → mein elements likhte hai

function in list

max(list) → maximum element deta hai [Print Karta hai]

min(list) → minimum element

sorted(list) → sorted copy data hai

list.append(element) → Element ko last mein insert Karta hai

list.clear() - Delete Karta hai sare elements ko

list.count(element) → kitne _{baar} elements hai list mein

list.remove(element) → Remove Karta hai element ko

list.reverse() → Reverse Karta hai list ko

list.pop → last element delete Karta hai

Sets - collection of unique element in no order

(iska order nhi hota)



Dictionaries - Key value pairs hoti hai python mein

Tuples - Ordered collection of items.

immutable

Loops

FOR loop -

→ Keyword declare the loop.

for i in Range(1, 100): f colon

Print i

indented.

condition for iteration/ Repetition.

while loop - yeh to tab stop hota hai jab condition false hoti hai.

i = 1

while i < 5:

keyword

Print(i * i)

i = i + 1

indented

Break - used to come out of loop.

i = 1

while i < 5:

Print(i * i)

i = i + 1

if (i == 3):

break

if 3 pr break hogi statement
baaki time chalegi.

continue - used to skip rest of the loop & transfer the controls to the top of the loop.

i = 0

while i < 5:

wapis
yaha
ajayega

i = i + 1

if (i == 3):

continue

Print(i * i)

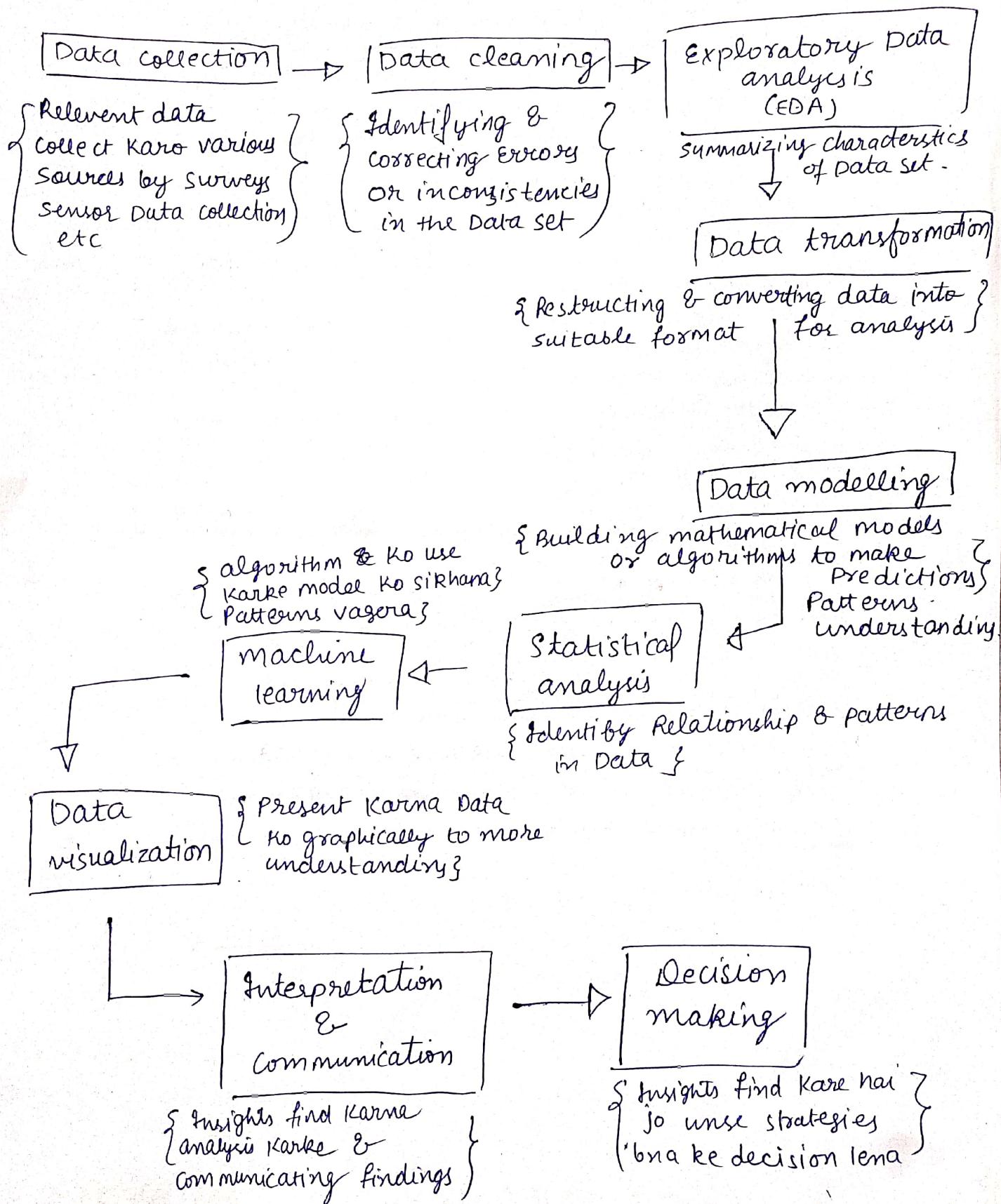
i = 3 ke lie

wali statement skip hogi

Data wrangling - process of removing errors & combining data sets to make them more accessible & easier to analyze

Data analysis

Process of inspecting, cleaning, transforming & Modelling data with the goal discovering useful information, drawing conclusions & supporting decision making.



Data importing -

↳ process jisme external files / data ko system mein dalte hai or analysis, manipulation, interpretation Karte hai.

↳ fundamental process hota hai

↳ Done by python libraries like "Pandas".

#

```
import pandas  
csv_file = Pandas.read_csv(file name)  
type(csv_file).
```

(work for
1D info
series)

data frame
(work for 2D
info)

If file is Excel file

```
excel_file = Pandas.read_excel(filename)
```

If file is JSON file

```
JSON_file = Pandas.read_json(filename)
```

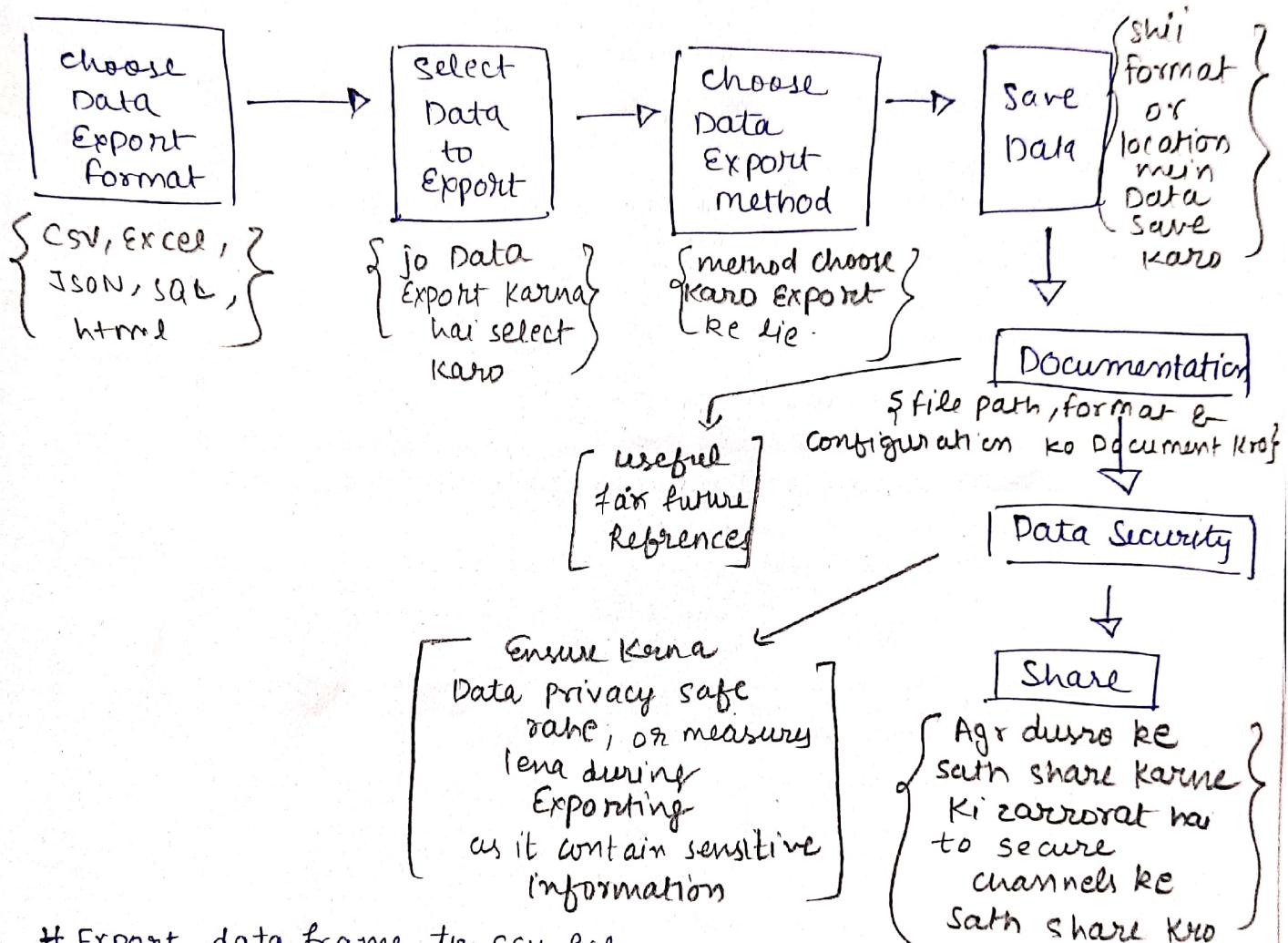
If file is html file

```
html_table = Pandas.read_html(url)
```

Data Exporting

Process jisme save or write karde hai data system mein or use external file or storage medium mein store karde hai.

Process



Export data frame to csv file

```
CSV-file.to_csv(file name)
```

Export data frame to Excel file

```
Excel-file.to_excel(file name)
```

Export data frame to json file object

```
json-file.to_json(file name)
```

Export data frame to pickle object

```
pickle-frame.to_pickle(file name)
```

Export data frame to html table

```
html-table.to_html(filename)
```

Basic Insights from Data set -

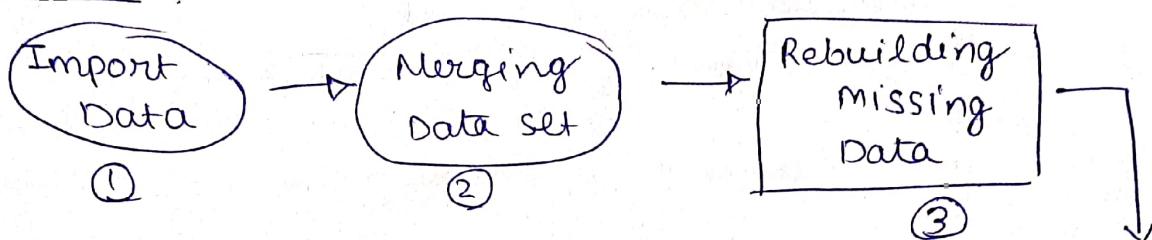
Extract karna Basic insights ka data se involve karta hai
Performing Exploratory data analysis jo understand
Karta hai characteristics, patterns or trends ko data set mein.

- ① Descriptive statistics - like mean, median, mode, range, variance.
Standard deviation offer karte hai Quick understanding of central tendencies.
- ② correlations - Identify karna Relationship between variables
pta lagana ki vo kaise ek dusre ko affect karte hai.
- ③ Trends & patterns - visualization representation (line chart, histogram / scatter plot) se Revel kar skte hai trend or patterns pta lagat skte hai for future.
- ④ Group comparisons - Data set ko split karna er Basic criteria pe jo differences or similarities show karega groups ke beeh.
- ⑤ Predictive modeling - machine learning algorithm predict karte hai outcome or future trends State hai Based on past
- ⑥ Anomalies or outliers - unusual Data points ko detect karna .
takri or deviate karna significantly tak
yeh impact na kare prediction ko.
- ⑦ Missing Data analysis - patterns / reasons ko samajna
missing values ke peeche ka
- ⑧ time series analysis - Examine karna ki raise variable change ho rke hai time ke sath, identify karna Seasonality or long term trends.

Data cleaning

↪ Data mein se Duplicate, incorrect, corrupted, incorrectly formatted info ko hatana is Data cleaning

Process

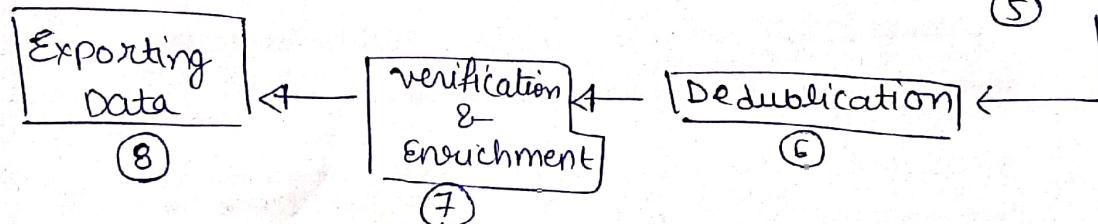


Standardization

④

Normalization

⑤



Standardization - Converting format to only one particular for all

Normalization - Sari values ko common scale pe karna adjust karna scale ko.

Deduplication - Extra value ya copies jo hoti hai Datee ki unko hatana

Verification & Enrichment -

↪ Missing values jo file ki hai & jo bhi changes kie hai unhe verify karna or store kar kar Raphna

Data cleaning

Pandas offer karne hai data frames (jo built-in hota hai)
Data ko clean Karne ke lie.

Renaming variable.

```
Product_df = product_df.rename(columns = {'cost': 'Price',  
                                         'shippingLocation': 'destination'})
```

Insert new column.

```
Product_df['NumberofProduct'] = 1
```

Checking for NULL values

```
product_df.isnull()
```

Boolean Data frame return karega, agr value null mogi to true(0)
varna false(1)

Kisi particular ko check Karana ho (column Particular ko)

Ki is column mein koi value null hai ya nahi

```
Product_df['Price'].isnull()
```

Dropping Rows & columns with Null values

Drop all rows having any Null value

```
Product_df.dropna(axis=0)
```

Drop all columns having any null value

```
Product_df.dropna(axis=1)
```

Drop all rows having less than k null values

```
Product_df.dropna(axis=0, thresh=5)
```

Drop all columns having all values

```
Product_df.dropna(axis=0, how="all")
```

Replace values -

Replace all occurrence of 0 & -1.

Product_df.replace(0, -1)

Filling NA values with some values.

Fill all NA values,

Product_df.fillna(0)

Grouping Data

Grouping on data by seller name

Product_df.groupby('seller name')

Apply statistics on grouped data -

Getting mean of cost associated from each seller.

Product_df['Price'].groupby(Product_df['Seller Name'].mean())
mean niklega

counting values

Product_df['Price'].groupby(Product_df['Seller Name'].count())
count karega value

Minimum value of cost of each seller.

Product_df['Price'].groupby(Product_df['Seller Name'].min())
minimum value

Maximum value of cost of each seller

Product_df['Price'].groupby(Product_df['Seller Name'].max())
maximum value of

Median value of cost of each seller -

Product_df['Price'].groupby(Product_df['Seller Name'].median())
Median statya

Std of cost .

Product_df['Price'].groupby(Product_df['Seller Name'].std())
standard deviation

Missing Data -

- ↳ common problem in machine learning.
- ↳ any information jo absent ho ya unavailable ho dataset mein.
- ↳ It occurs ~~where~~ when there is gap or Empty field in collected information.
- ↳ can happen due to human error during data entry.
 - ① or Malfunctioning Sensors in data collection devices.
 - ② not choosing any answers during surveys,
 - ③ or information was not recorded

Deal Karna missing Data ke sath Kafi crucial hai as yeh affect karta hai accuracy, Reliability ya kisi bhi conclusion ko jo data set se nikala jata hai.

Methods to handle missing values

(1) Remove features which contains missing values -

Simplest way is to analyze which feature contains the missing value & delete them, model ko baaki features (like jo ~~columns~~ bohot columns bache hain), unke sath train kia ja skta hai.

Ye Sif tab use ho skta hai jab huma Data set kaafi acha ho or missing values boht kam features ki ho, or unka correlation jo ho dependent variables ke sath vo high na ho.

(2) Remove Rows contain missing values -

hum Rows ko delete kar kte hai missing values ko hatane ke lie pr pher agr hamare test data mein

vo Rows present hogi to machine learning algo fail ho jayegi correct output predict karne mein,

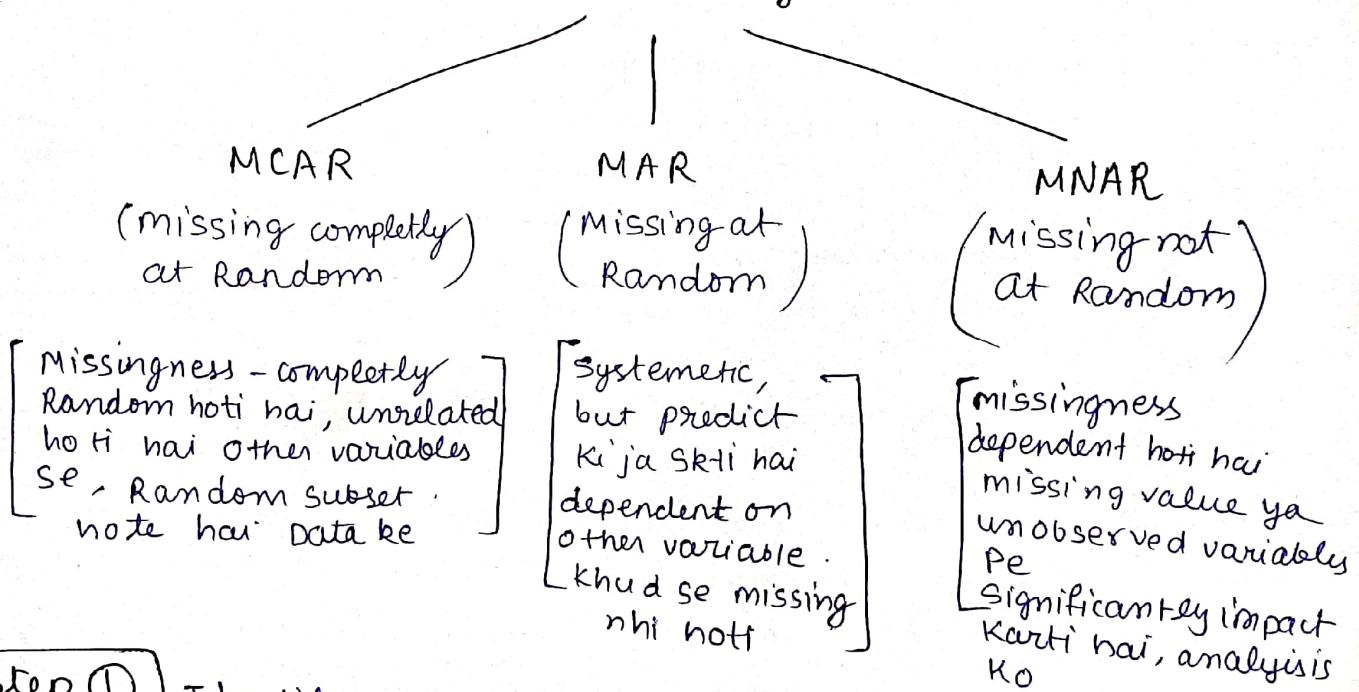
or age large no. of Rows mein missing data hogi to vo training data ko reduce kar dega.

③ Replace missing values with significant data

- ↪ well technique.
- ↪ values ko Replace kar dete hai different value ke sath Numerical values mein Replace kar sakte hai.
- Mean ke sath feature value ki.
- median ke sath feature value ki.
- zero ke sath.

Categorical ke lie attribute a ek new category "default"

Types of missing values



Step ①

Identify Karo missing value ko kis type ki hai
mehr numerical / categorical (NAN / NULL)

Step ②

Understand Karo nature ko kya hai (MCAR/MAR/MNAR)

Step ③

- (i) Kam hai to Row/Column ko delete kardo.
- (ii) Imputation (fill Karo values if numerical)
 - if categorical - default.
 - if numerical $\begin{cases} \text{Mean} & \text{if numerical} \\ \text{Median} & \text{if categorical} \end{cases}$

Step 4

jisse dataset ko test hoga usme update Karo values

Impact - Impact dekho jo hamari missing values fill ki
hai, humne unka impact dekho kya pade,
Prediction pe kya padega

Documentation - Save Karo ki kisme values ko / Riski values
change ki thi as if Baad mein dekhni padi to

Quality / completeness - sawal missing values handle hogi
or data set ki quality check karni hei

UNIT-IV

Mathematical & scientific applications for Data analysis

Mathematical

{ Network analysis
Relationships }

Graph theory

{ Numerical analysis }

Computational Mathematics

Differential Equation

{ used to
model systems
applied to
time series
data }

Mathematical

Probability

{ deals with
uncertainty &
chances of
happening or
non happening }

{ Numerical analysis }

Statistics

{ central tendency - mean, median,
mode, variance describe
main features of Data set }

Linear algebra

[transforming, scaling
data,

for dimensionality
Reduction - singular value
decomposition]

Population, sample

Regression, correlation

Correlation

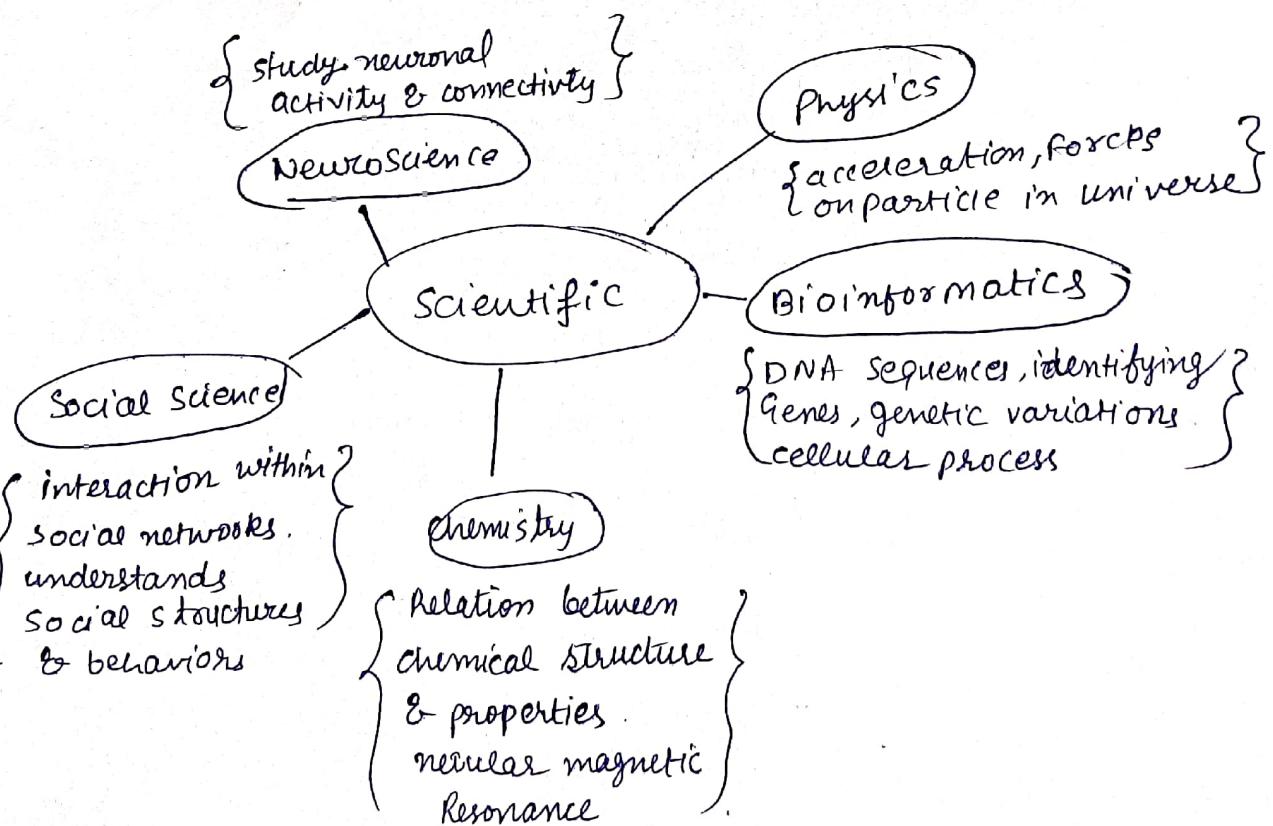
Mean, median,

Mode, Variance

Describe

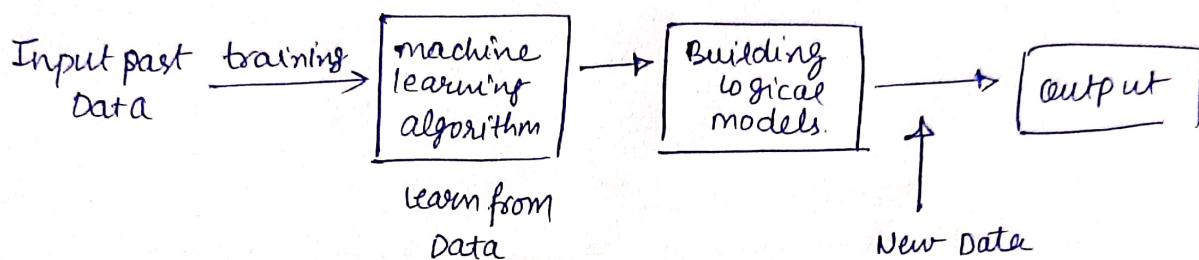
Main features

Data set



Machine learning -

- Subset of artificial Intelligence
- focuses primarily on creation of algorithms, so computers to enable them to independently learn based on the data set / previous experiences etc.
- without being explicitly programmed, machine learning enables a machine automatically learn from data, improve performances from experiences & predict things.



Features of machine learning -

- uses data to detect various pattern in given set.
- learn from past data & improve automatically.
- Data driven technology.
- Deals with huge amount of data

Need for machine learning

- Rapid increment in production of data.
- Solve complex problems, difficult for a human.
- Decision making in various sector including finance.
- Find hidden patterns, & extract useful info from data.

Classification

- 1. Supervised Learning
- 2. Unsupervised Learning
- 3. Reinforcement learning

Applications

- Image Recognition, Medical Diagnosis
- automatic language translation, self driving cars

Supervised learning

Type of machine learning hota hai, jaha algorithm learns Karta hai labelled training data se, correct output data milta hai.

meaning: Input data ka corresponding output labels or target values hoti hai.

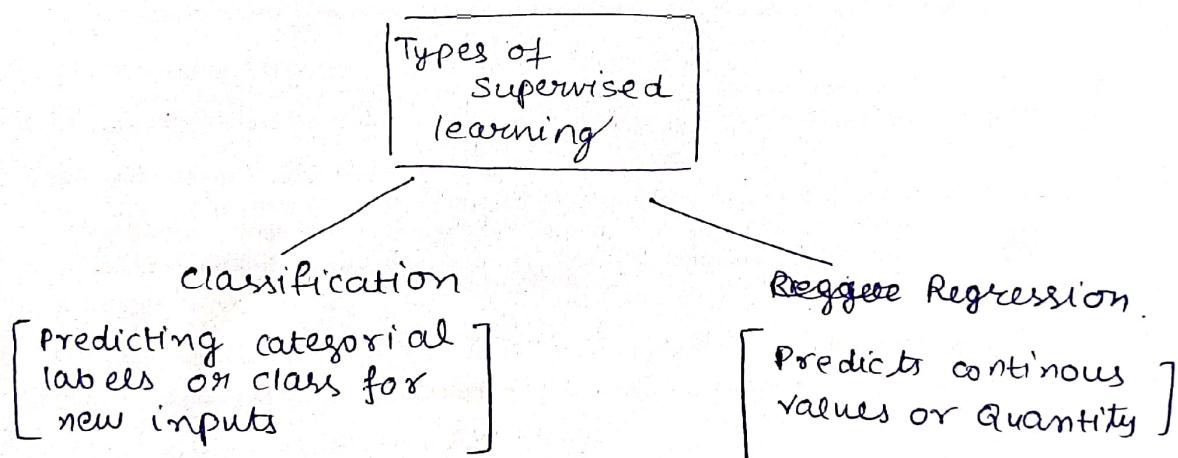
→ Supervised learning ko goal hota hai mapping function Se learn karna (input se),

[Application] - fraud detection, spam filtering

[Labelled data] - consists of input-output pairs.

[Learning algorithm] - algorithm learns from labelled data (Patterns & Relationships)

[Prediction] - Model can predict outputs



Workflow

Data collection → Preprocessing → Feature Extraction



Prediction ← Evaluation ← Model training

Machine learning

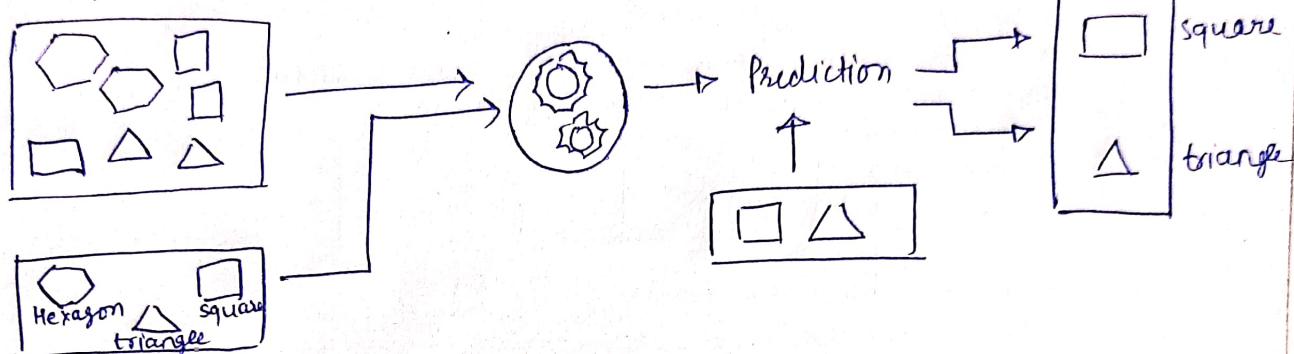
- Sub set of artificial intelligence
- Focus primarily on the creation of algorithms that enable a computer to independently learn from data & previous experiences

Supervised learning -

- ↳ type of machine learning jisme machines ko trained karte hai labelled data ke saath, or correct output predict karte hai.
- ↳ Aim : To find a mapping function to map input variable (X) and output variable (y)

↳ applications - Risk assessment, Fraud detection, Spam filtering.

Working -



steps

- 1) Determine type of training Data set .
- 2) collect / gather the labelled training data .
- 3) split the training Data set into training dataset , test dataset & validation Data set .
- 4) Determine input feature of training data set .
- 5) determine suitable algorithm .
- 6) Execute algorithm .
- 7) Evaluate the accuracy of model by providing the test set
age model ne shi output predict kia , iska matlab model
shi hai

advantages

- model can predict the output on the basis of prior experience .
- Exact idea about the classes of object .
- used to solve real world problems .

Disadvantages

- not suitable for handling complex tasks .
- Training Required lots of computation time .
- Enough knowledge about classes of object .

Unsupervised learning -

type of machine learning where algorithm is trained on unlabelled data & tasked with finding pattern or relationship, Structure within Data Bina Explicit Guidance ke

→ Goal ~~hota hai~~ explore karna or understand karna kise ander ke structure ko or usse extract karna useful insights bina labelled diagram ke.

Types

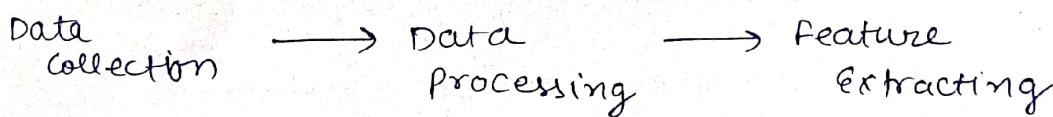
Clustering

Group similar data points together Based on some similarity or dissimilarity

Dimensionality Reduction

Reduce no. of features while Retaining important info.

Workflow -



Evaluation ← Model training

Advantages

- used for more complex tasks as compared to supervised learning.
- easy to get unlabelled data

Disadvantages

- result might be less accurate.
- more difficult

Supervised

- trained using labelled data
- predicts the output
- Produces an accurate result
- needs supervision to train the model along with output
- takes direct feedback to check if predicting output correct or not.
- Supervised
 - Classification
 - Regression
- Data is provided along with the output

unsupervised

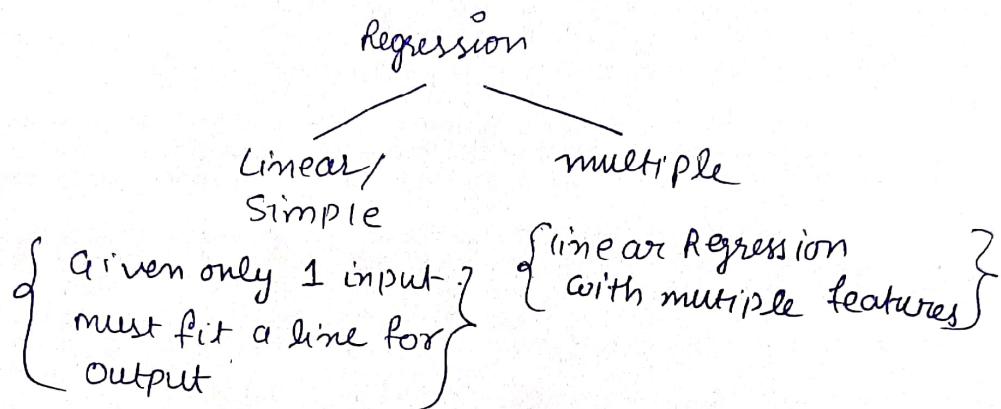
- trained using unlabelled data
 - find the hidden pattern
 - may give less accurate result as compared.
 - doesn't need any supervision to train the model
 - doesn't take any ~~any~~ feedback
- unsupervised
- clustering
 - Dimensionality Reduction
- only input Data is provided to the model

Regression (technique) (supervised learning)

→ Predicting the continuous valued output from input.

Example application of Regression

- Predicting stock price
- Forecast sales of moth.



Classification - (supervised learning technique).

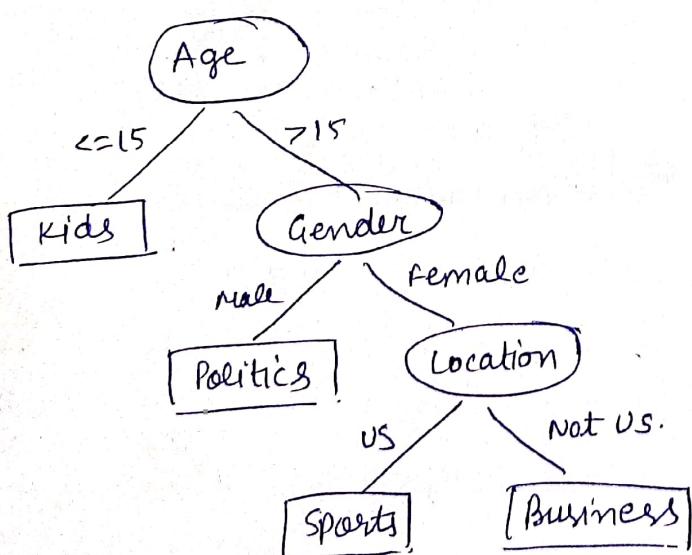
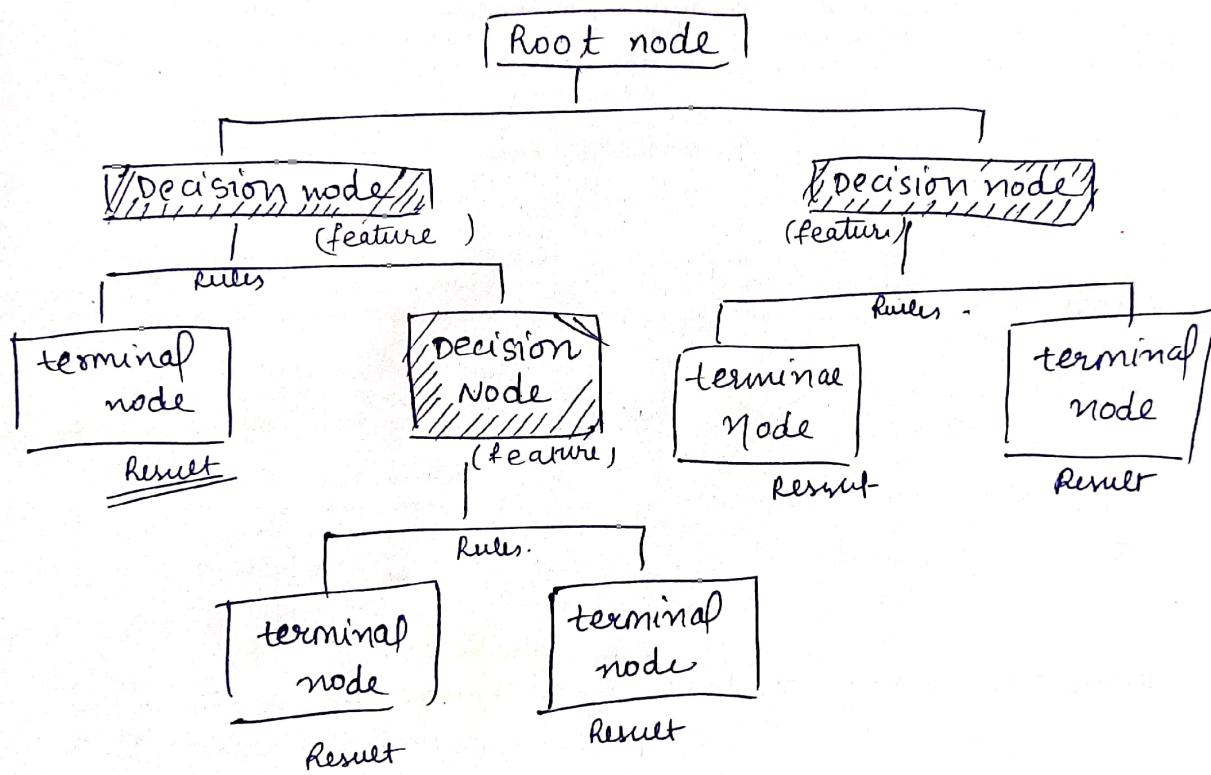
- with the help of features & categorical output, one has to predict the output.
- models are classifier.
 - Role to learn the relationship between input variable & categorical outcome.
 - classes.

Ex classifying emails - 'Spam hai kinh'.

classify document - from which genre a thing belongs like - In newspaper it is related to sports, politics, Education, market etc.

Decision tree -

- technique of classification.
- hierarchical tree structure which consists of a root node, branches, internal & leaf nodes.
- internal nodes of tree denotes feature, branches denotes - the Rules.
- leaf nodes denotes the Result of algorithm.



Advantages of Decision tree

- (1) capture non linear relationship between the Data
- (2) Easy to Explain to other . (doesn't need any complex mathematics knowledge to understand the result)
- (3) can handle both numerical & categorical .
- (4) fast process .
- (5) handle irrelevant attributes .
- (6) not based on assumptions .
- (7) cleaning on data not needed .

Simple tree

↳ Early stopping - Hum tree ko Banana start Karte hai, and stop Kar dete hai ek depth pr agar

tree ko Bnana tab stop Karte hai ya to jab.

(1) validation error or decrease na hora ho.

(2) ya pher jab boht kam data Bacha ho kise node ke pass, usse majority ko leke prediction kerte hai.

↳ Pruning - Build Karte hai pehle tree ko, or uske Baad tree ki branches ko cut Karna start Kar dete hai better generalization error ke lie.

Decision tree algorithm - operate Karta hai data ko

Analyze karke.

Recursive algorithm., It Repeats until Reach a leaf node.

Step 1. Start Karo tree ko Root node ke sath, jo complete dataset contain larka hai.

Step 2. Best attribute find Karo data set mein using attribute selection measure.

Step 3. Divide Karo Data set ko subset mein, contain larka value according to best attribute.

Step 4. Decision tree node generate Karo.

Step 5. Recursively new decision tree Banao, or use Karo @ Subset data set ke jo step 3 mein bnae the.

or is process ko continue larka jae tek ap data ko or classify karne ke kabil na bacho., or un nodes ko final node / leaf node boldo.

Disadvantages

- (1) Kahin na kahin decision tree ko stop hona chalta hai,
Varna high error milne ke chances hote hain.
- (2) faces problem of overfitting.

Decision tree

classification

Dependent variable -
Categorical

Statements ka use nahi
hota

Regression

Dependent variable - continuous.
like Ex - Salary Based age,
years of experience.

Regression tree mein, kisi thi
node pe mean value di jati
hai sab values ki, as
prediction.

Decision tree learning

- algorithm used to learn the decision tree from training data
- work Karta hai top to bottom
- best attribute check Karta hai kisi thi point pe, or learn
Karta hai, create Karta hai decision node uske lie.

Decision Boundaries :-

Gini score -

split of different class attribute at any node
Best value = 0 (when node divide data perfectly)
worst value = 1 (when both branches have classes in same proportion)

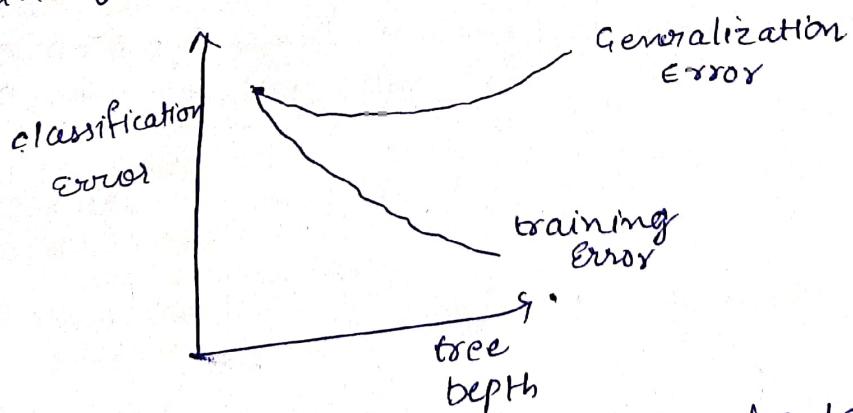
Goal - create a model jo reduce kare Errors ko.

Error - defined as no. of incorrect predictions divided by total no. of examples taken.

Decision tree problem

face karta hai problem of overfitting,
agar hamare pass boht sare features hai to hum bs.
split karke branches hi bñate rahenge, or jab end mein
boht kam data bachega har grp mein to vo generalized
solution hi hoga.

isise Training error to decrease hoga w.ith increase in
depth of tree but generalization error increase hoga



simple tree better hai overfitting ke term mein
isme choose kia jata hai depth or complexity ko
decision tree ki use validation ke Based Pe.
age, simple & complex dono tree's same validation
error denge / approximate denge, to simple ko
choose karzenge.

Random forest - (Ensemble technique)

↳ forest of trees.

↳ contains large no. of decision tree jo help Karti hai
decision lena mein.

↳ tree same strategy se bnae jata hai joaise decision
tree bnae jate the or decide kia jata hai majority
votes se.

Random forest advantages over Decision tree

→ Missing data ko handle kar leti hai

→ Regression or classification dono ho jate hai

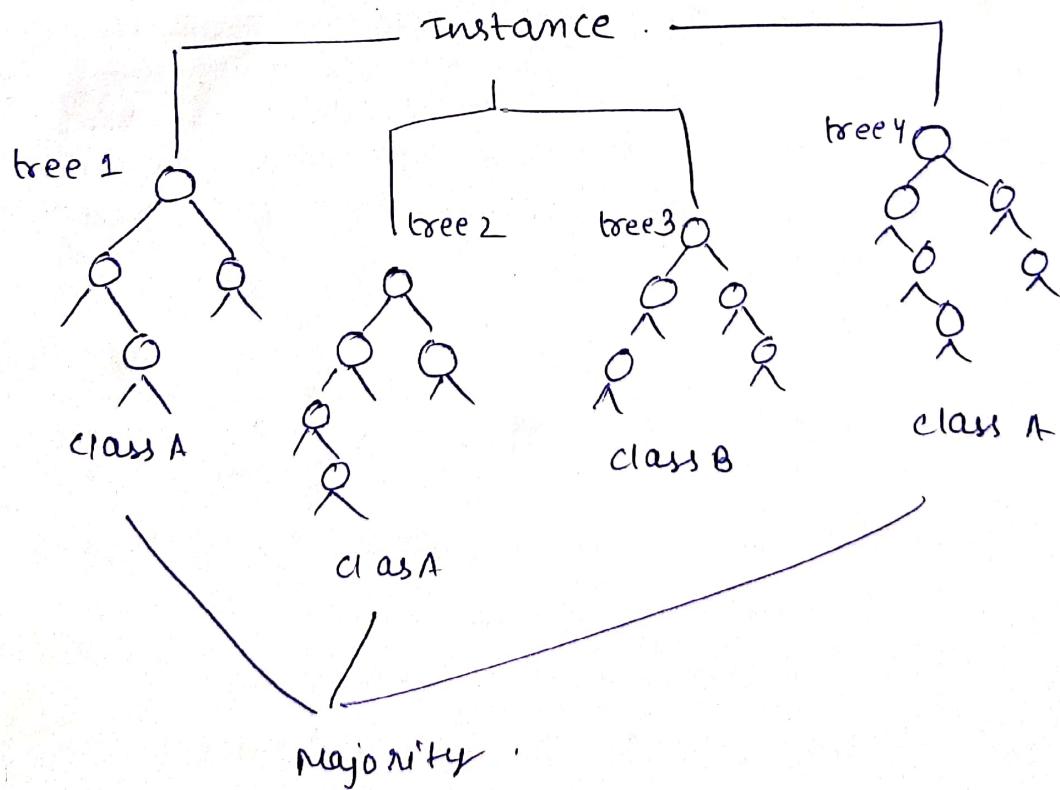
→ overfitting nahi hoti

Disadvantages -

→ model ke upr kaafi kam control hota hai

→ complex hai

→ difficult to Explain.



Naïve Bayes

- ↳ highly used in text classification.
- ↳ advantage - Speed ↑, fast & make prediction very easily with high dimension of data.

- ↳ probabilistic classifier.

- ↳ naive because yeh man ke cheta hai ek feature dusre feature ke independent hota hai.

~~Explain~~

- ↳ Bayes theorem se apply hota hai

$$\boxed{\text{Posterior Prob.}} \quad \boxed{P\left(\frac{x}{y}\right)} = \frac{P(y) \underbrace{P(x)}_{\substack{\rightarrow \text{Likelihood} \\ \text{Evidence}}} \rightarrow \text{Prior}}{P(y)}$$

Clustering - way to group the objects in a way that

- similar object single group mein Rakte-
 - dissimilar alg alg
- ⇒ intercluster similarity ↑
intracluster dissimilarity ↑

Feature	Predictive	Trends
<u>objective</u>	future outcome predict karta hai based on historical Data .	Identify karta hai patterns & tendencies over time .
<u>Nature</u>	future focused	Past focused .
<u>Data usage</u>	historical data build karke prediction karta hai .	Examine karta hai historical data identifying karta hai patterns .
<u>techniques</u>	machine learning methods	time series analysis
<u>Decision making</u>	Aims to inform decision making for future events .	Assist in making decision Based on historical Data .
<u>Example</u>	weather forecasting .	seasonal variations

Predictive mining

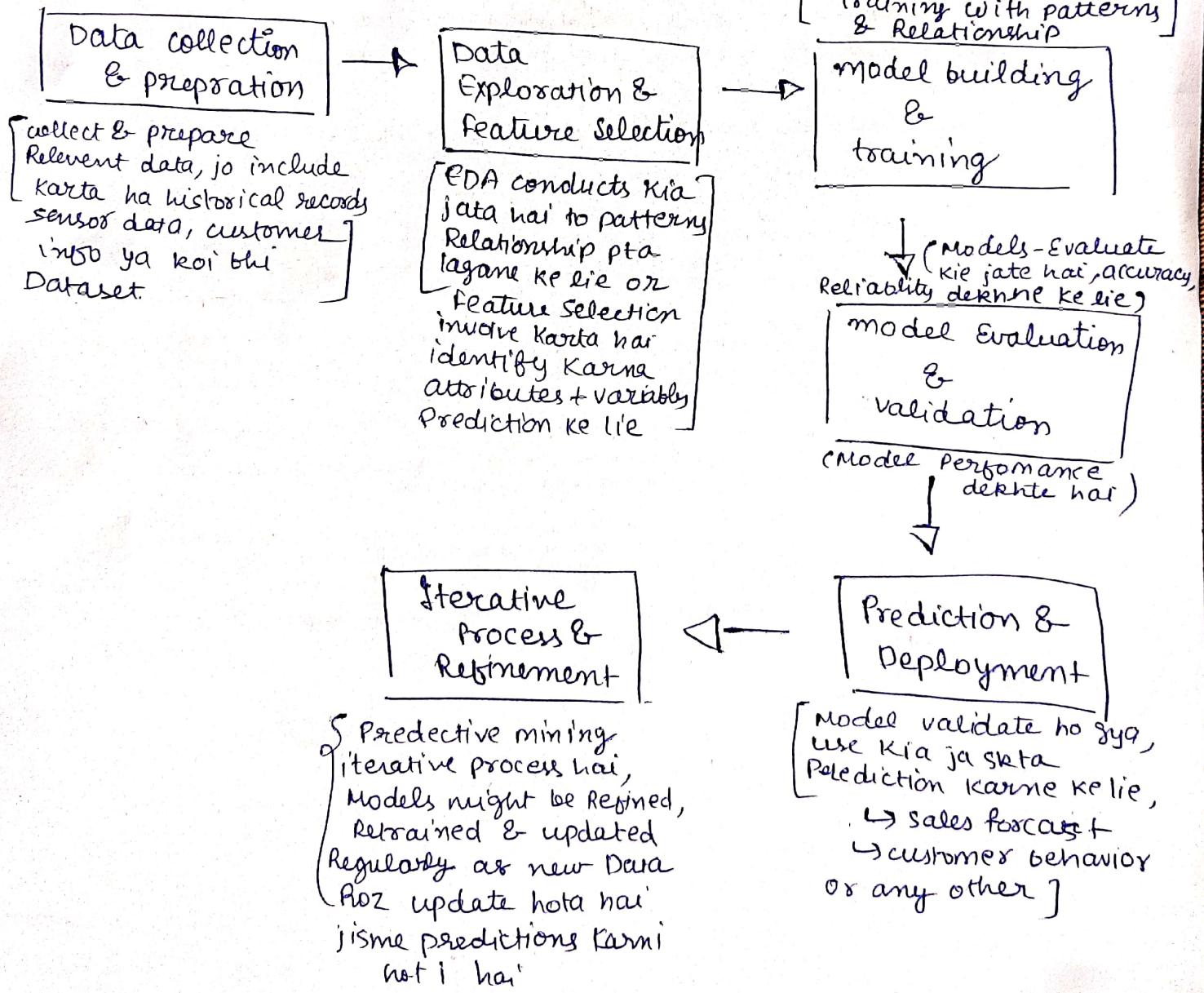
Prediction + Data mining

Practice of using Data, Statistical algorithm and machine learning techniques, Pattern identification & predict future outcomes based on historical data

Exploration & analysis of large Data sets to discover meaningful patterns, Relationship or insights, understanding/ Perception

machine learning algo, Regression, decision tree neural network, model training with patterns & Relationship

Process



Trend analysis

Trend analysis -

Examining data over time jisse identify karate hai
consistent patterns or tendencies ko.

↳ used in finance, Economics, marketing, technologies.
to make informed decision Based on historical data.

↳ It Reveals insights to direction & behavior of variable
or set of variable ek specified period ke lie.

Key aspects -

1. Data collection :- Relevant data collect karta hai over a specific time period., measurable variable se hota hai.
2. time series - Data ko chronologically arrange le carte hai,
time series create karne ke lie. Each data point series ka represent karta hai variable ke value, Ek specific time ke lie.
3. visualization - plotting graphs jisse trends easily find out kar sake.
4. Identification of patterns - data ko analyze karenge or recurring patterns nikalenge.
5. statistical technique - Regression analysis, mean, median, mode & all.
6. forecasting - planning & decision making, trend identify hone ke Baad.
7. factors influencing trends - dhundo un factors ko jo trends ko affect karte hai.
8. long term/ short term trends -
kitne time tk ka trend hai.

Recommendation system -

Applications or platforms jo Recommend karte hai useful things to the user according to interest.

Goal of Recommendation system -

- # Present Karna items user ko jisme vo interested hai.
- # Yeh selection or increase Karne mein help karta hai sales ki

Recommendation Engine -

1. Popularity Based.
2. Content Based
3. Classification Based.
4. Collaborative Filtering.
5. Model Based.

Popularity Based -

→ Based on usage of items by all users

(jo item sbse jyade use kia jata hai, sbse popular manajata h)

islie vo sbko suggest kia jata hai

→ Based hota hai past usage count pe.

→ yeh recommendation system used hota hai news article
(jisme Recommend kia jata hai item new user ko)

→ iska drawback hai, yeh assume karta hai Har user same type ka hai, or sbko same recommendation deta h.
without even knowing about the profile geography or other attributes.

Important

content Based Recommendation Engine -

- Based hota hai feature pe item ke,
- items ko compare karte hai or recommend karte hai,
jiske features mil rhe hote hai us item se jisko user ne
higher Rate kia hota hai.
- Memory Based System hota hai kunki it memorizes the
items.
- but inme training data large hota hai or fit nhi kar pata
memory mein islie prefer nahi karate

IMP

Collaborative filtering -

- ↳ Based on the idea of similar interest
- ↳ it depends on the choices of two users

collaborative

user-user

Based on similar user, on their item purchase

User ek hi hoga pr uski choices dekhni jayengi Kya-kya item usko pasand aa rhe hai

item-item

→ Based on similar items purchased by users

User do honge but unme similarity dekhni jayegi, ki

Kya kya item dono ko Pasand aa rhe hai,
ya ek ko arha nai to }
dusre user ko vhi recommend }
hoga

{ agar ek item pasand aa
rha hai to usi re jaise
item recommend honge }