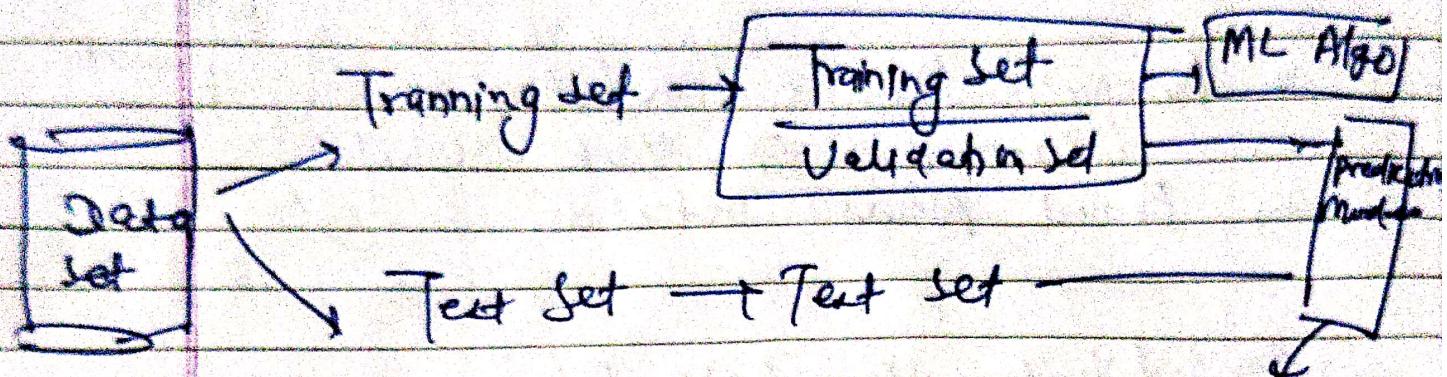


11

Need of Data Set ↗



Top Sources

→ Kaggle Data Set

→ Amazon Data set

→ UCI ML Repository → Classified for go
Bktl Already

→ Google dtsr Search
Engine

Clean data
most dtsr set Clean
E.g.

→ Microsoft Dctg Set

→ Awesome public Data Set Collection
(GitHub repo)

→ Sckit - Learn Data Set

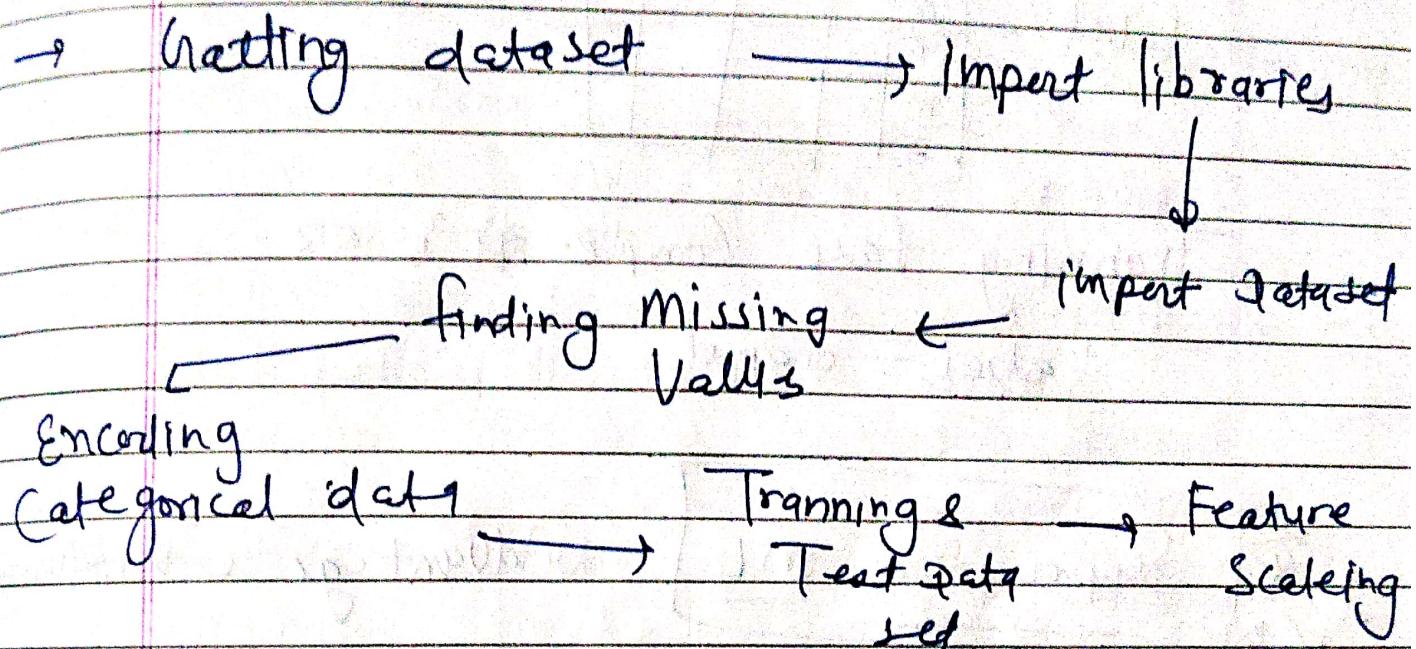
Toy
(practise) → scikitlearn datasets
→ computer vision dataset (VGG16)
host dataset
designing & training
data gen



Scanned with OKEN Scanner

12 (2) Data Preprocessing in Machine Learning

It is a process to convert data at suitable form जिसका है



Depended on Independent Variable

$$x+y = z$$

Google Colab file → Data Preprocessing .ipynb

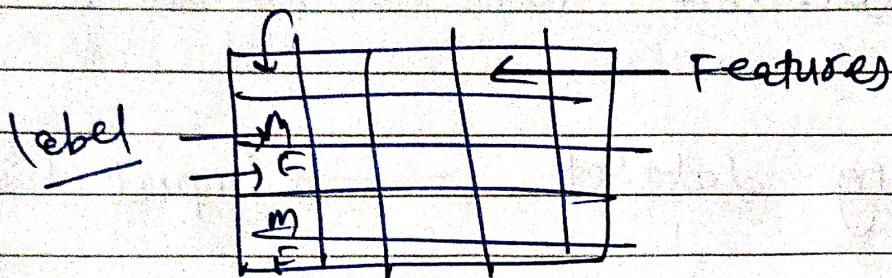
All Step we are doing step by step

(13)

Features and label in ML

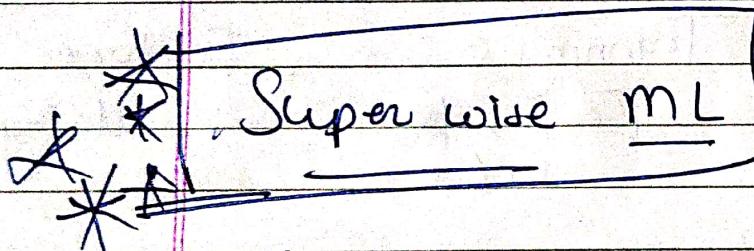
Notetaking
Date _____
Page _____

Column → Features



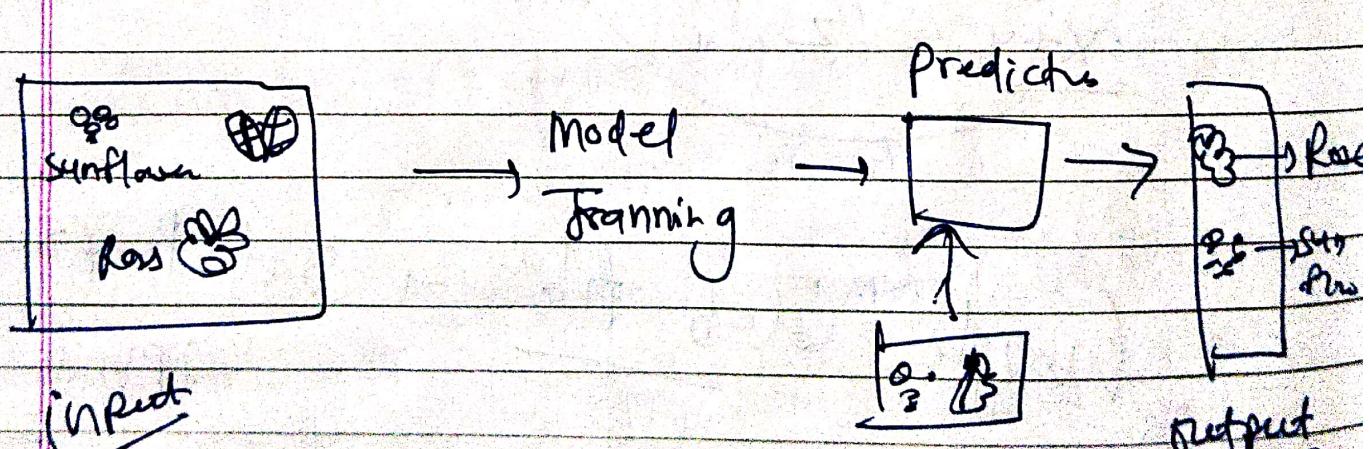
Training fase (Complexity के बारे में)

(label = output)



Advantage & Disadvantage

→ Input Variable की match करते ही फिर किए



2 types



Scanned with OKEN Scanner

Regression

get output ~~not~~ Numerical

~~not~~

Ex: 1, 0, 22, 23

Unique Values

→ ~~use~~ use ~~in~~ ~~for~~ ~~it's~~ ~~it's~~
→ ~~it's~~ depended ~~on~~ ~~it's~~
independend ~~it's~~ ~~it's~~
relation ~~it's~~ ~~it's~~ ~~it's~~

→ Continuous Values

Ex: Weather forecast
market Trends

Classification

get output ~~not~~

Categorical ~~not~~

True false Yes, No

Male female

Day, night

→ Most popular Algo.

Random Forest

Decision Trees

Logistic Regression

Support Vector Machine

Most popular Algorithms.

① → Linear Regression

② → Polynomial Regression

③ → Regression Trees

Advantage of Supervised Learning

- ① Full control over ML debug model
- ② easily fast and determine the number of classes
- ③ we can

DisAdvantage of SLF

- Have limited scope
- Label data set is expensive and time-consuming
- Wrong prediction
- इसमें एक supervisor होना है जो कि निर्देश में दिया है
- supervisor ML पर भी जिस data से Train करते हैं वह label data देता है
- 2 types
 - Regression (Numerical)
 - Classification (Categorical)

Ex:- Spam filter
 → label की target करता है
 जैसे win lottery, win, Jackpot
 mail में जित की spam की जल देता

→ प्राप्त हो Training data से अलग
 जैसे है, हमें वो > 100 result

Unsupervised Machine Learning

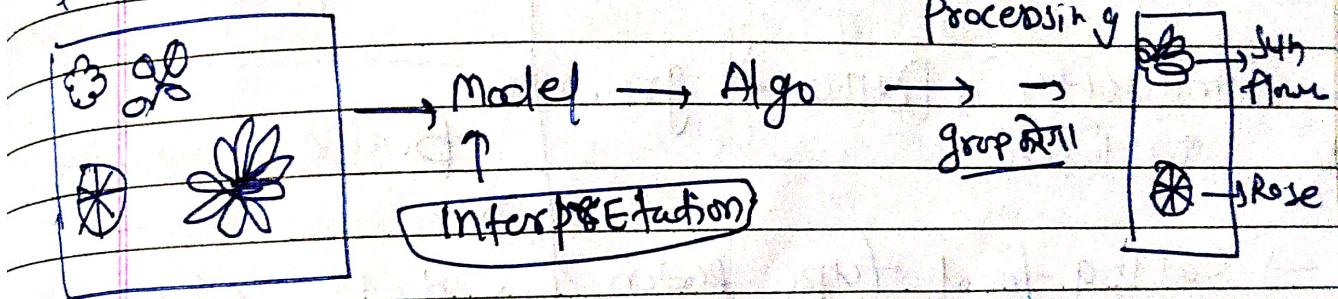
Unlabeled data

Model को find करे hidden patterns and insights से उनकी data को रहे हैं

Data को sort करा similarity के लिए।

→ नहीं previous knowledge को लेता है

Dataset



Unlabeled data

→ यह model pattern search करेगा और Category के Base ac बताएंगे कि flower मिसाएंगे

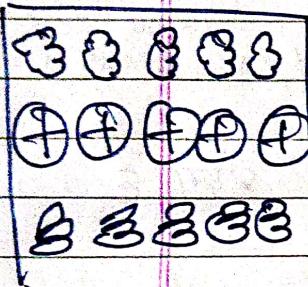
The goal of Unsupervised Learning is to group unlabeled data according to the similarities, patterns and difference without any prior training of data.

Type of Unsupervised ML

Clustering

Association

Raw Data



\rightarrow Algo \rightarrow

Output



Brade + Butter + Egg =

Egg + Br + Choco =

Choco + Br + Chipp =

Chipp + Br + Choco =

Most frequent

similarities के आधार पर group
किया जाता है।

Brade + Butter

\rightarrow eating food type fusion

\rightarrow Onion, Tomato & mix के

उनसे अलग² नहीं होता है

of all objects में
relationship देखेगा।

Ex- Glossary example

most people

Home Laptops हमें नहीं लगते हैं और वरीदत है

इससे Marketing effective
की गति है।

Brade के बाद ही कोई हो
Butter की ओर भी हो।

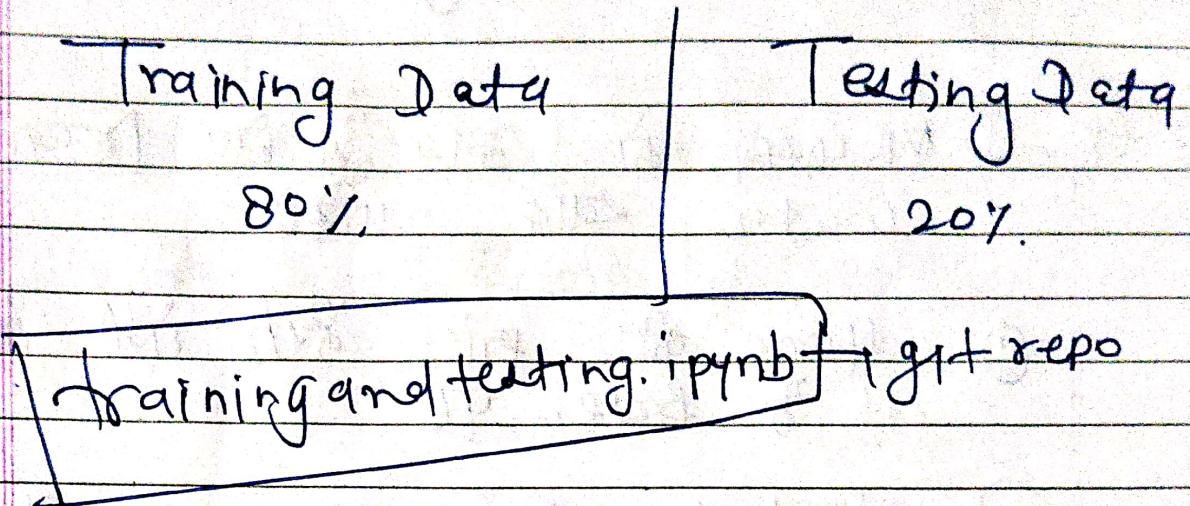
Unsupervised Learning Algorithms

- K-Means Clustering
- KNN (K-Nearest Neighbour)
- Hierarchical Clustering
- Neural Networks / Deep learning
- Single Value Decomposition
- Distribution Models
- Principal Component Analysis
- Apriori Algorithms

~~Advantage~~

- ① Most Complex task to perform
 - ② It is help in finding pattern in data
 - ③ saves lot of manual work and expense
- in early stages

⑪ * Training & Testing Data in ML



⑫ Linear Regression single variable

It is part of supervised ML

→ It's famous and easy

→ इसमें हम जो Value predict करते हैं
वो Continuous ही वाले or
numerical Value ही रखते

Ex Cost age, salary temp.

Price 2 types

① Linear Regression with single Variable
② Linear Regression with multiple Variable