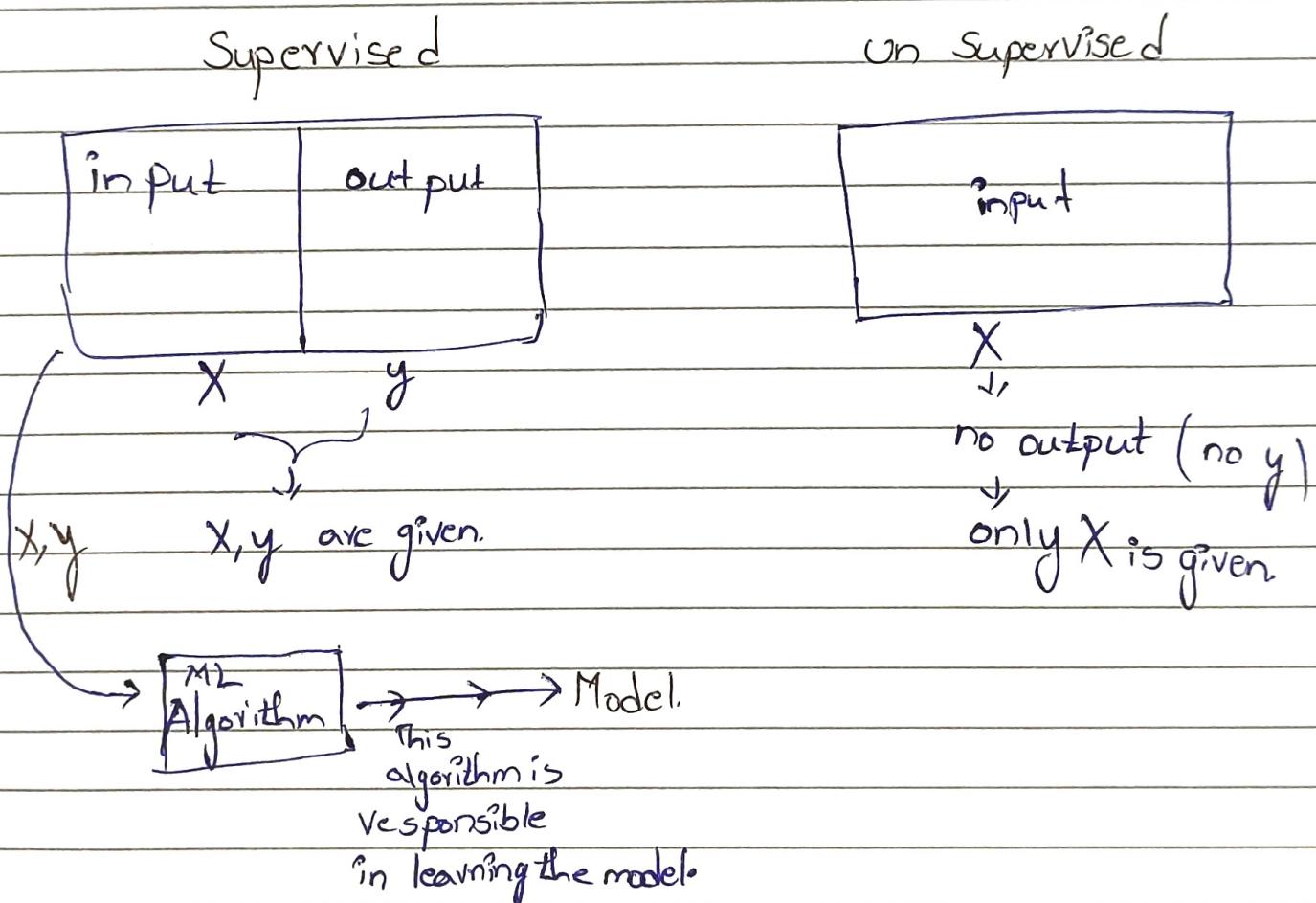


Q) How to identify whether a problem is unsupervised learning?

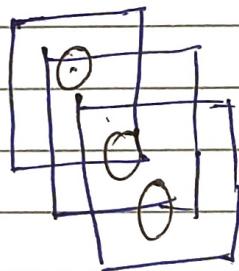
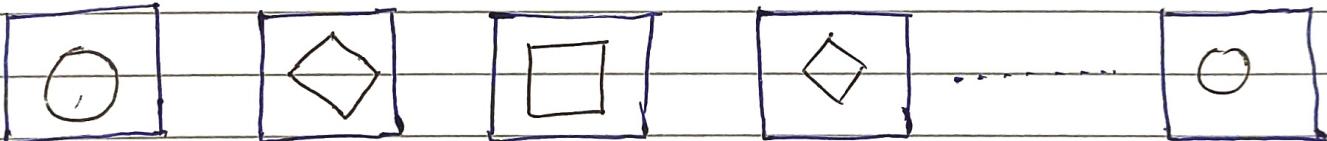
A) There is no output label associated with each input



$$X \rightarrow \boxed{\text{Model}} \rightarrow \hat{y}$$

L) This model will predict Future/unseen data.

- In classification, the output column will tell us whether we need to perform classification / Regression.
- If output column is not there then we need to go with unsupervised learning.
- I am having 1000 cards which consists of different objects
So, we will group similar cards.



- ↳ There is nothing like prediction in unsupervised learning.
we are trying to recognize the pattern & we will group it together.

Data mining →



Collecting the data & recognizing the pattern

Ex:- Extract gold (in mining) → Looking for gold

↳ Similarly extract some useful information from the data.

↳ Looking for pattern

Dataset Definition →

$$D = \{ (x_i, y_i) \mid x_i \in \mathbb{R}^d, y_i \in \mathbb{R} \}$$

↓ ↓
Input output
↓ ↓
Input target "d" no. of columns

→ \mathbb{R} → Real Value feature

Dataset



Continuous Value



Regression

$$D = \{ (x_i, y_i) \mid x_i \in \mathbb{R}^d, y_i \in \{-1, +1\} \}$$

↓ ↓
Input output
↓ ↓
Input target "d" no. of columns.

Input is having
"d" no. of
columns.

y_i is having either
-1 (or) +1
↓
Classification.

$D_n = \{(x_i) \mid x_i \in \mathbb{R}^d\} \rightarrow$ For un Supervised learning.

→ In data set each row is a vector

→ So, here x_i is a vector

→ x_i is having 'd' no. of components ("d" dimensional like 2-D, 3-D, ...)

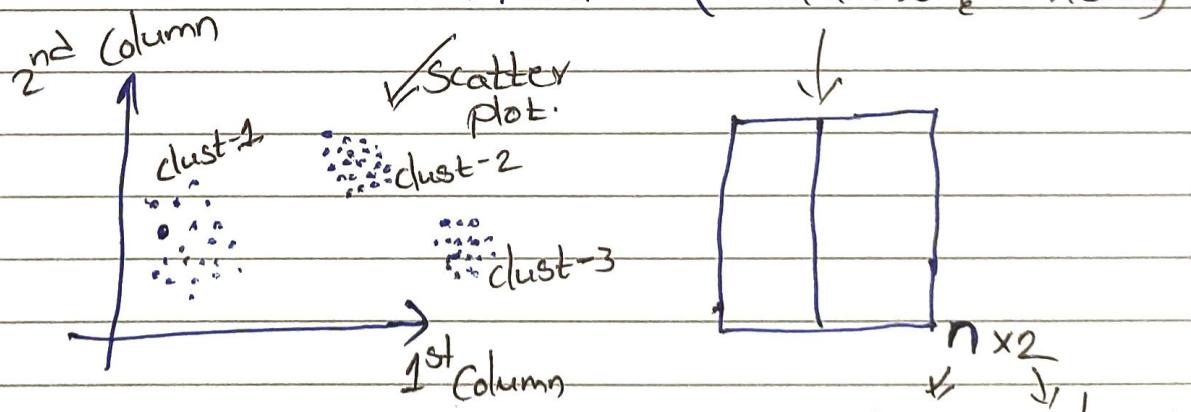
→ x_i is Real Value

ex:- I want to store a single value in 'a' then
 $a = 5$

→ If I want to store two values in 'a' then

$a = \begin{bmatrix} -1 \\ 2 \end{bmatrix} \rightarrow$ Vector \rightarrow ITS 2-D

$\hookrightarrow x_i \in \mathbb{R}^2$ (2 columns & n rows)



Observation \Rightarrow we have 3 Groups / cluster

Cluster: → Group of Similar data points.
↓
distance

- * points in a cluster should be cluster close together
↳ (Intracluster distance → small)
- * points in two different clusters should be far as possible
otherwise two clusters will combine.
↳ Intercluster distance → high

2. 2 3
1,2 → Similar
1,3 → dissimilar

I can say it using distance

Clustering is UnSupervised learning

Intra → within

~~Intra~~ → Between
Inter

K-means, K-means++, hierarchical, DBSCAN.
clustering

classification

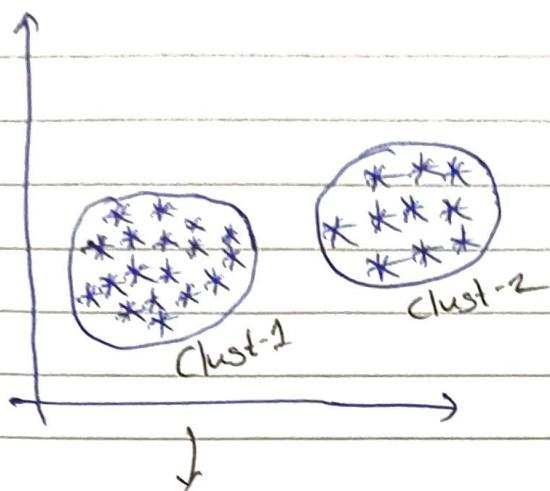
↳ Task → separate different classes.

Regression

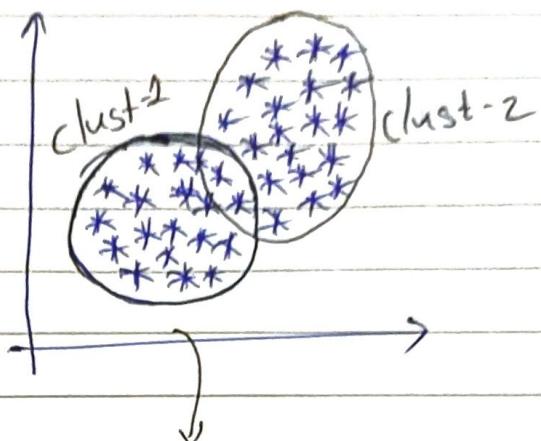
↳ Task → try to fit Best fits the data.

Q) what if, if a customer / data point belongs to all the clusters.

Ans: It means clusters are overlapping clusters are not properly defined.



clusters are properly defined



clusters are not properly defined. (It is overlapping)
In this
↓ area

We need to do more Feature Engineering

So that this clusters stops overlapping each other.

In Real world → We will see many times that clusters are overlapping

In that case there are many other algorithms we can look at GMM

↓ ↗ probabilistic approach → Gaussian Mixture Model approach

It will assume Normal Distribution of all the clusters

Applications of Unsupervised Learning

1. E-commerce

→ jio, airtel → some customers are not happy with the services provided by Airtel so they will move to other (Jio)

When we do unsupervised learning on top of this data then we know this group of customers are not happy.

So Rather than providing offers everybody, if they give discount to unhappy customers then it is good for them.

(Anyway happy customers will use Airtel (won't use Jio)
even though you give offers/Not)

That's why we used to get messages in whatsapp.

Stating 50% discount,....

2. Image Segmentation.

↳ Grouping of similar pixels

↳ after grouping, we typically apply ML techniques to perform object detection.

↳ Grouping of similar pixels

↓ in terms of ML

Group Nearest data points

↳ Vector.

(Self Driving Car)

↳ It has object Detection)

↓
It came into picture in recent years

→ But this Segmentation (self driving car) has been

used in 1980's by US Defense. But not come into this world because of Computational power.

→ Because these Vectors are not having need high Computational power.

↳ Now only laptops, computers are having high

Computational power (like SSD). and it was available
to everyone easily.

→ That's why now it is getting popular.