

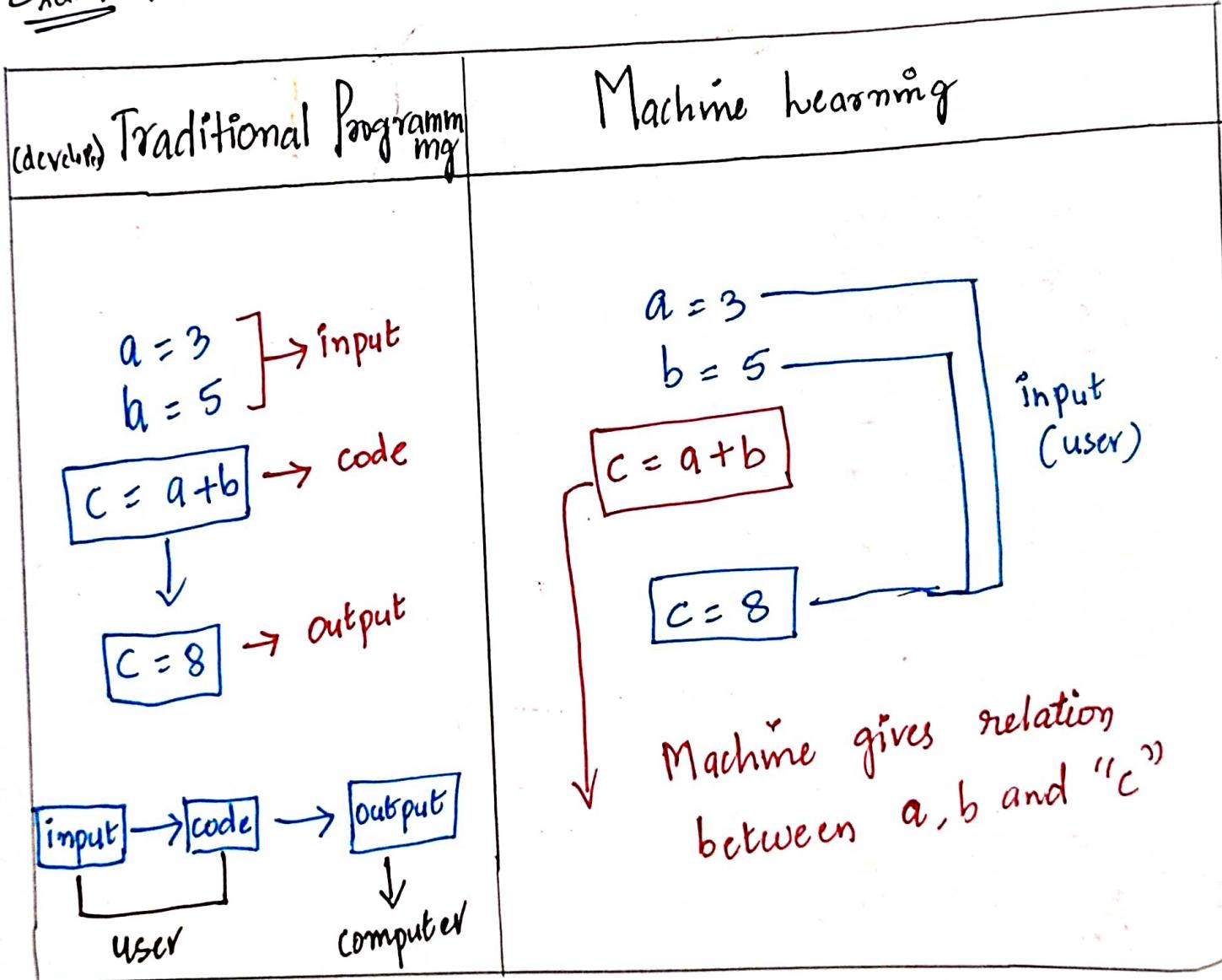
# → Machine learning

## Introduction

Ques:- What is Machine Learning?

Ans :- Machine learning is The subField of computer Science that gives "Computer Ability to learn without being Explicitly programmed"

Example:-



input	output
X	Y
1	10
2	20
3	30
4	40
5	50
6	?

$\rightarrow g$  (Prediction)  
 $"y = 10x"$

input have more columns

$x_1$	$x_2$	$y$

$$y = f[x_1, x_2]$$

$$ax_1 + bx_2 = y$$

## MODELLING

Q. How I can say it is 6?

Ans :- \* Identifying The between "x" and "y"

$$\therefore y = 10x \Rightarrow \text{MODEL}$$

\* Substitute "6" in "x" value

$$y = 10[6] = 60$$

$$\therefore y = 60.$$

## PREDICTION

## Machine learning

\* Machine  
 \* Identifying

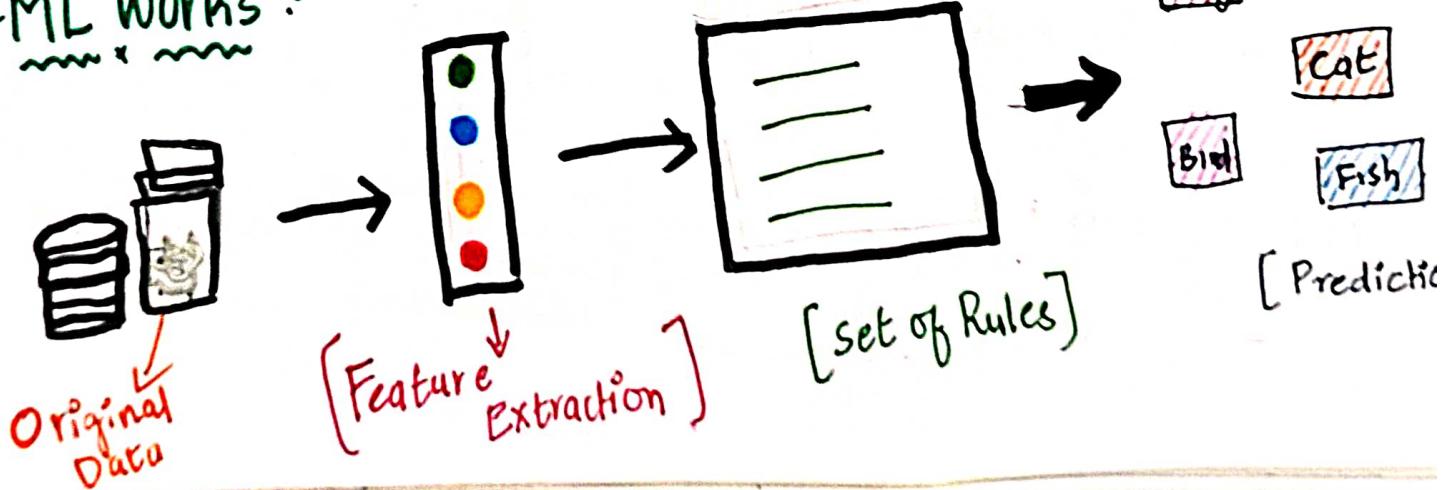
Output

learn The given Data.

The relation between input and  
 is called Machine Learning.

\* Prediction

\* ML WORKS :-



# ⇒ Major Machine Learning Techniques

- \* **Regression / Estimation**
  - Predicting **Continuous Values**
- \* **Classification**
  - Predicting **Discrete Values** / **Item class** / **category of a case**
- \* **clustering**
  - (No prediction)  
- Finding the **structure of data**; **Summarization**
- \* **Associations**
  - Associating **Frequent Co-occurring items / Events**  
Ex: Buying Brush/paste.
- \* **Anomaly detection**
  - Discovering **abnormal** **Unusual Cases**
- \* **Sequence Mining**
  - Predicting next Events; **Click stream** (**Markov Model, HMM**)

\* Dimensional Reduction

- Reducing The size of data [PCA]

\* Recommendation systems

- Recommending Items.

Ques What is Supervised Learning?

Ans:- We "Teach The model" then with that knowledge,  
it can Predict Unknown (or) Future instances.

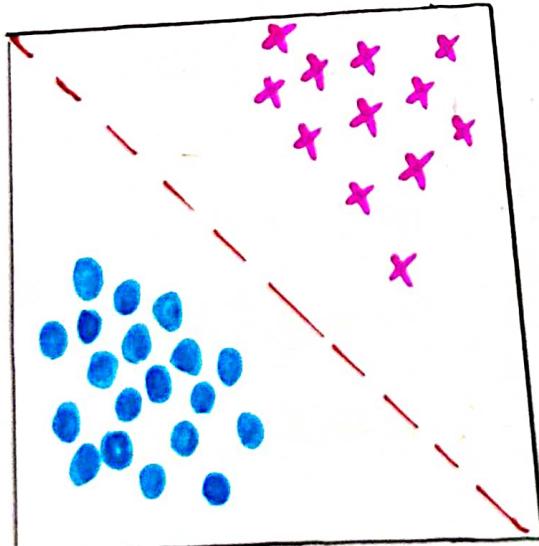
→ Where Ever The Output Variable is Available  
is called as "Labeled data".

"In Supervised Learning Model, "The Algorithm Learns on a  
Label dataset , To generate Reasonable predictions For The

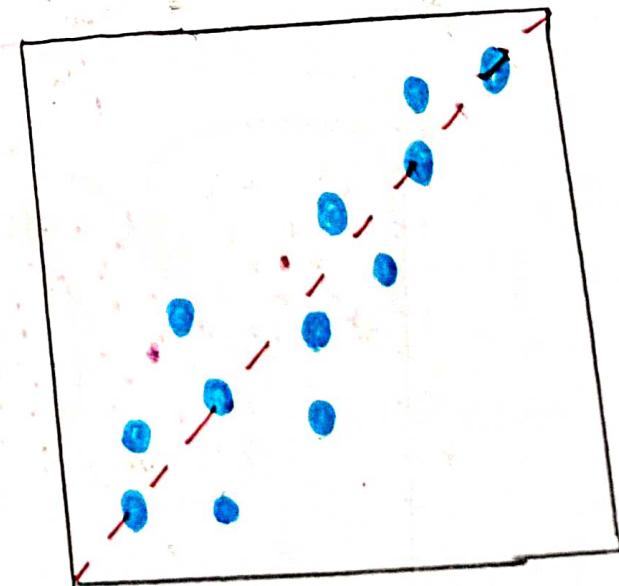
response to new data.

1. Regression [Continuous Data]

2. Classification [Discrete Data]



\* Classification



\* Regression

Ques :- What is UnSupervised Learning?

Ans :- UnSupervised Learning, we have

Unlabeled Data

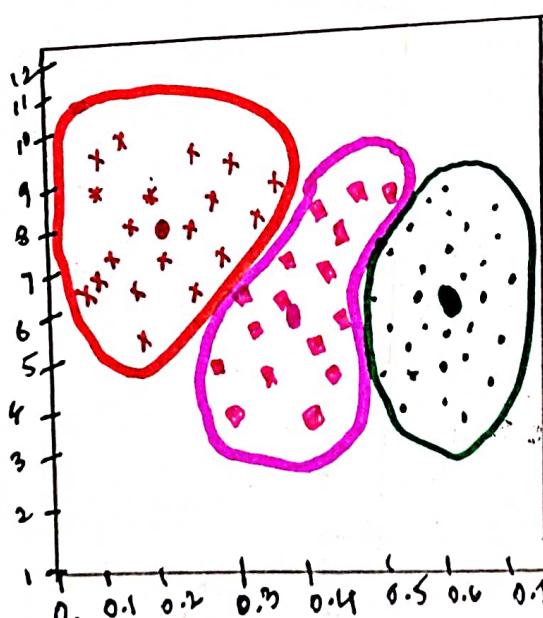
i.e. **No Output Feature Available.**

"**Unlabeled Data**" that the algorithm tries to make sense of by **Extracting Features**, **Co-occurrence** and **Underlying patterns on its own.**

1. clustering
2. Anomaly Detection
3. Association
4. Auto Encoders

Ques :- what is clustering?

"clustering" is **grouping of data points** (or) **Objects that are somehow similar by**



\* clustering

## \* Supervised Learning

- \* Regression
  - classifies labeled data
- \* Classification
  - Predicts trends using previous Labeled data
- \* Has more Evaluation methods than Unsupervised Learning.
- \* Controlled Environment

## \* Unsupervised Learning

- \* Clustering
  - finds pattern and grouping from Unlabeled data.
- \* Has Fewer Evaluation Methods than Supervised Learning
- \* Less Controlled Environment

Dt: 6/4/22  
10:30pm

6/4/22  
3:30 pm

## Sklearn

↓  
Scikit Learn

### # Feature Scaling

From Sklearn.preprocessing import StandardScaler

df[ "Age-sc" ] = sc.fit\_transform( df[ [ "Age" ] ] )

Using Pandas  
input, we giving as Data Frame

Out :- We get DataFrame.

### # missing values

From sklearn.impute

df[ "CLMAGE" ] = pd.DataFrame( mean\_imputer.fit\_transform( df[ [ "CLMAGE" ] ] ) )

import SimpleImputer  
using Pandas.

Out :- We get array.

Data Frame  
Input, we giving same

Example :-

### User Defined Functions

```
def add(a,b)
```

$$c = a + b$$

```
return float(c)
```

→ add(3,4)

Out: 7.0

→ Float value

```
def add(a,b)
```

$$c = a + b$$

```
return c
```

→ add(3,4)

Out: 7

```
def mains(a,b)
```

$$c = a + b$$

$$d = a * b$$

$$e = a \% b$$

$$f = a - b$$

```
print("sum": c)
```

```
print("mult": d)
```

```
print("div": e)
```

```
print("sub": f)
```

→ Saving in "Jack.py"

From Jack import maths

↓  
module

↓  
Function

\* Because, function is created in a such a way That  
answer should be float. / But Both The Functions are same  
↓ input

Only output is different.

\* Same,  
Simple Imputer, Standard Scaler and More Functions

From sk.learn are InBuilt Function.

\* We Can't Remember Every function. So, We Take The

Help function and see The Output Variable Type.

\* If Out is Array. We change (or) Convert into  
[Pd. DataFrame] before The

Data Frame, Just writing

Function.

→ Estimator API

## Modelling

```
# import 'ML' Algorithm  
# model = MLAlgo()  
# model.fit(x-train, y-train)
```

## Prediction

```
# yPred = model.predict(x-test)
```

## Evaluation

```
# Comparison (y pred, y test)
```

# Accuracy

Ex:-

X	Y
1	10
2	20
3	30
4	40
5	50
6	?

# model = Algorithm [ $y = 10x$ ]

yPred # model.predict [x-test]

$$y = 10x \quad [x=6]$$

$$y = 10 \times 6$$

$$\text{output} = 60$$

Relation :  $y = 10(x)$



ML. Algorithm [model]

Students



Learn



write The Exam



Evaluate



grade / percentage

## \* MODEL Evaluation \* → Train-test-Split

\* Train and Test on Same Dataset

\* Train / Test Split

EX:-

Case 1

\* In a class Sir given + Interview question

Train → 100 Q & A



Test → 30 Q

80% marks (%)

You're a good student, But  
Only for this class Only.

\* If whatever The Data I have  
given [100 Q/A]. if we ask  
30 Q/A from the same dataset.  
it gives better performance, But.  
Only for that dataset it  
works, not for any new  
dataset (Q/A)

Case 2

100 Q/A

70  
[Training]

30  
[Testing]

↓  
For These 70 (Q/A) Given To Machine, not  
to understand, But given Only For  
Understand Concepts Only. not For reading data

Now, We ask new 30 (Q/A), not from  
70 (Q/A) that we have given, and  
Evaluate Accuracy For Training and For  
Testing, Separately.

## Example

	input						y (Output)
	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$y$
0	2.0	3.0	3	8	4	3	196
1	8.0	4.8	2	22	11	9	121
2	3.0	3.2	1	44	22	16	136
3	4.8	3.8	8	2	118	22	225
4	9.0	5.2	4	3	22	11	224
5	6.8	4.4	3	44	114	18	230
6	8.0	3.2	8	8	55	16	255
7	2.8	8.9	2	3	88	14	264
8	8.0	10.10	5	2	63	12	288

If This is input, what is out (Predict by model)  
Even Though we have "y" value we consider it, But Machine has Predict same Value.

Original Question

Ans

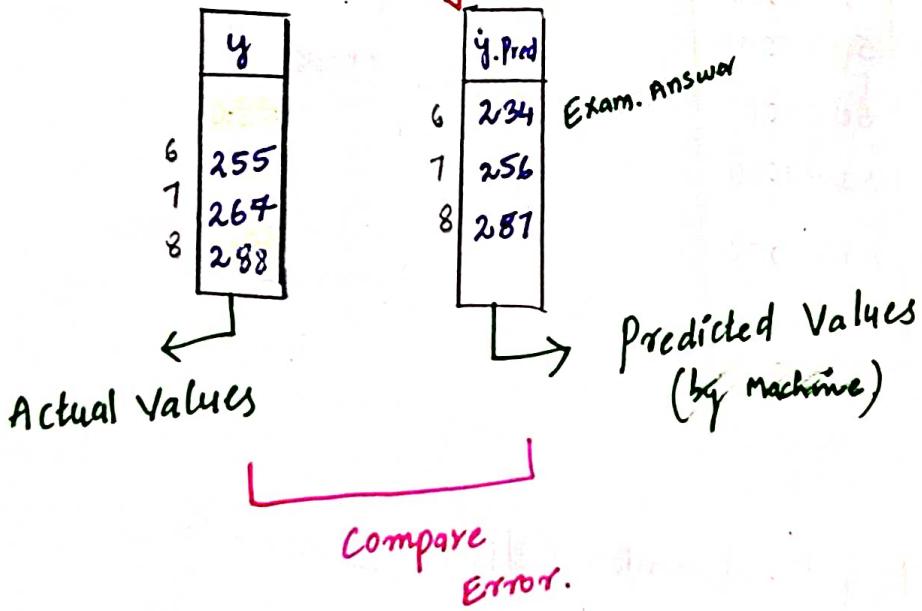
\* In Testing, we give input data that we are separated First.

and check with Accuracy or "y" value (output)

\* In Testing, we don't take "y" variable (output) because To

Check Given data, Model is Working (or) not.

MODEL



Code :-

```
# Import numpy as np
# Import pandas as pd
# Import matplotlib.pyplot as plt
%matplotlib inline
```

```
# df = pd.read_csv ("homeprices.csv")
```

df

**Out** :-

	town	area	Price
0	Chennai	2600	5500000
1	Chennai	3000	5650000
2	Chennai	3200	6100000
3	Chennai	3600	6800000
4	Bangalore	2800	5850000
5	Bangalore	3300	6150000
6	Bangalore	2600	6500000
7	Bangalore	3600	7100000
8	Hyderabad	2600	5750000
9	Hyderabad	2900	6000000
10	Hyderabad	3100	6200000
11	Hyderabad	3600	6950000

# Here we use same dataset  
of Encoding. and we divided  
them into train / test.  
it will easy to understand  
by taking same Data set for  
further steps.

# creating dummies

```
# df_dum = pd.get_dummies (df, dropfirst=True)
```

```
# df_dum
```

**Out**

Now, this  
Changed  
dataset, ←  
Is we Consider  
Further to  
Train/test

Area	Price	town-chennai	town-Hyderabad
2600	5500000	1	0
3000	5650000	1	0
3200	6100000	1	0
3600	6800000	1	0
2600	5850000	0	0
2800	6150000	0	0
3300	7100000	0	0
3600	5750000	0	0
2600	6000000	0	1
2900	6500000	0	1
3100	6200000	0	1
3600	6950000	0	1

# Here Based on Area (sft) and town, we want predict The price

Independent Variable

# So... area, town-chennai, town-Hyderabad are → Input variable

# price is Output variable (y)

Dependent Variable

Because, it's ←  
Depending on "x" values To give output.

# Option To Consider x & y Variable

x = df\_dum.drop("Price", axis = "columns")

y = df\_dum.Price

# check "x" → Input Variables  
(or)

X  
Independent  
Variables

# Here we split the data into Train / test.

# stores in "x" and "y".

# check "y" → Output Variables  
(or)  
Dependent Variables

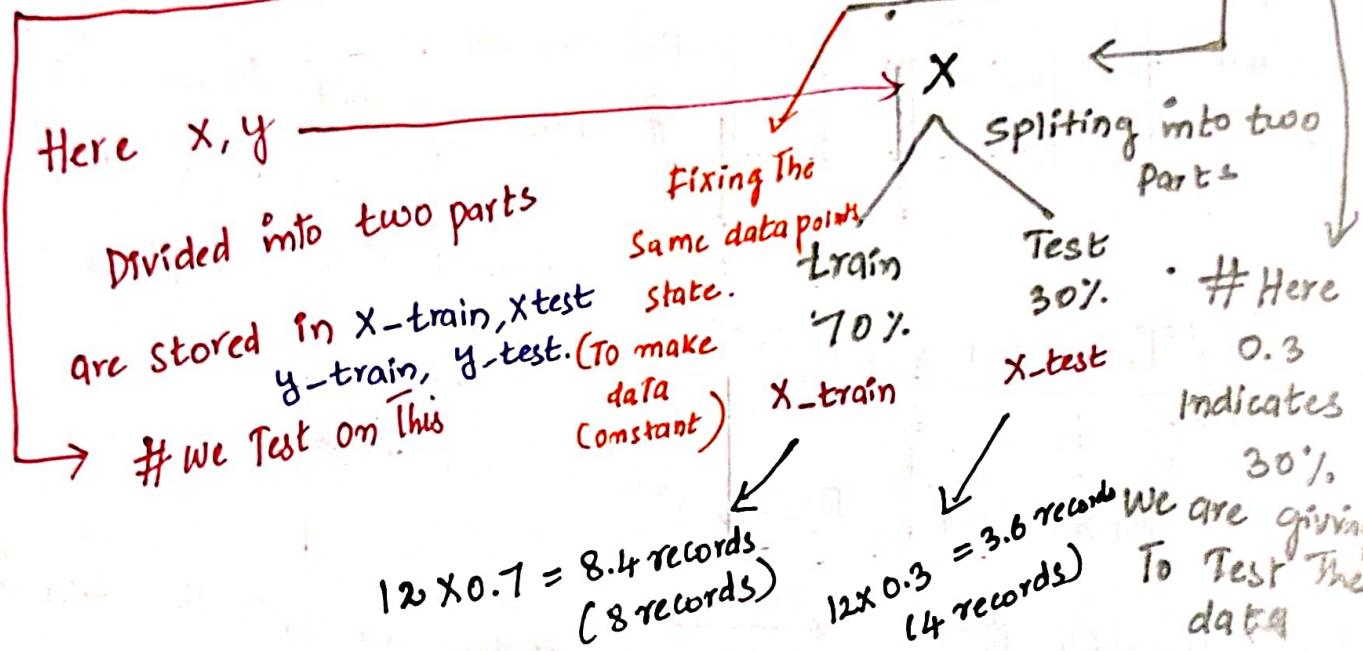
y

Out:	Area	town-chennai	town-Hyderabad
0	2600	1	0
1	3000	1	0
2	3200	1	0
3	2600	1	0
4	2800	0	0
5	3300	0	0
6	3600	0	0
7	2800	0	0
8	2600	0	1
9	2900	0	1
10	3100	0	1
11	3600	0	1

Out:	Price
0	5500000
1	5650000
2	6100000
3	6800000
4	5850000
5	6150000
6	7100000
7	5750000
8	6000000
9	6500000
10	6200000
11	6950000

from Sklearn.model\_selection import train\_test\_split  
 ↓                   ↓  
 Library      Module

#  $x\text{-train}, x\text{-test}, y\text{-train}, y\text{-test}$  = train-test-split ( $x, y, \text{test\_size}=0.3$ , Random\_state = 29)



"ML" Algorithm 1

\* Same Data

"ML" Algorithm 2

\* Same Data

# In order to compare which algorithm is working well and giving high accuracy.

# We have to same data for every algorithm and compare with all algorithm

code:

```
# X_train, Xtest, y_train, ytest = train_test_split(x, y, test_size=0.3,
Random_state = 29)
```

# X-train

Out :-

Area	town-chennai	town-hyderabad
7 3600	0	0
6 3300	0	0
1 3000	1	0
0 2600	1	0
8 2500	0	1
2 3200	1	0
3 3600	1	0
5 2800	0	0

# X-test

Out :-

area	town-chennai	town-hyderabad
9 29	0	1
4 26	0	0
10 31	0	1
11 36	0	1

# y-train

Out :-

Price
7 7100000
6 6500000
1 5650000
0 5500000
8 5750000
2 6100000
3 6800000
5 6150000

This Heading  
is not  
mentioned  
in my pc  
→ Only  
These  
Values.

# y-test

Out :-

Price
9 6000000
4 5850000
10 6200000
11 69500000

✓  
✓  
✓  
✓  
✓

4:00 AM

Dt: 7/4/22  
12:30 PM

### Step 3 :- Prediction

# Train

$$\text{MODEL} = 10x + 5$$

Substituting x values in -

$$\Rightarrow 10(1) + 5 = 15$$

$$\Rightarrow 10(5) + 5 = 55$$

$$\Rightarrow 10(9) + 5 = 95$$

$$\Rightarrow 10(36) + 5 = 365$$

$$\Rightarrow 10(12) + 5 = 125$$

$$\Rightarrow 10(4) + 5 = 45$$

$$\Rightarrow 10(14) + 5 = 145$$

Example :-

# Past Data

Input

	X	Y
0	1	10
1	5	15
2	9	36
3	14	45
4	12	12
5	4	88
6	78	56
7	2	44
8	14	68
9	36	52

# Based on  
The data  
model is  
formed (Equ)

# If we change  
some data (1, 2 records)  
our model also  
changes.

X	Y	y-Pred-train
1		15
5		55
9		95
36		365
12		125
4		45
14		145

# Test

X	Y	y-Pred-test
2		25
14		145
78		785

$$\Rightarrow 10(2) + 5 = 25$$

$$\Rightarrow 10(14) + 5 = 145$$

$$\Rightarrow 10(78) + 5 = 785$$

Step 1 :- Split The Data For Train/test

	X	Y
0	1	10
1	5	15
2	9	36
3	36	52
4	12	12
5	4	88
6	14	68

	X	Y
7	2	44
3	14	45
6	78	56

# Test

# Train

70%

Selected  
Randomly

30%

Step 4 : Evaluation

# Train	Y	y-Pred-train
	10	15
	15	55
	36	95
	52	365
	12	125
	88	45
	68	145

# difference b/w  
Actual "Y" - Predicted "Y"  
= ERROR

$$\begin{aligned} \text{Ex: } & 10 - 15 = -5 \\ & 15 - 55 = 40 \\ & 36 - 95 = 59 \end{aligned}$$

# Train\_Error

Step 2 : Modelling

$$\text{MODEL} = 10x + 5$$

# Test

Y	y-Pred-test
44	25
45	145
56	785

# Test\_Error

## \* Machine Learning Algo steps

### → MODELLING

1. Finding The relation between X & Y of given data

$x\text{-train}$ ,  $y\text{-train}$

# relation → "MODEL"

Ex:  $y = f(x)$   
 $y = ax^2 + bx + c$   
 $y = mx + c$   
 $y = 10x + 5$  (Equation)

### → PREDICTION

Substituting

the values of "x" in the given Equation and

MODEL calculating

"y" values.

### → EVALUATION

Calculating

the error between Predicted Values and

Actual Values

# Good Model

# Train-Error = 5% Accuracy = 95% ] [±5]

# Test-Error = 5% Accuracy = 95% ] ↓  
with difference

⇒ train accuracy ≈ Test accuracy  
with deviation of ±5%. # Good Model.

## # Overfitting

- train Accuracy = 90%  
→ Train Error = 10%.
- test Accuracy = 60%  
→ Test Error = 40%.

[Model] [Performs good in training]

But,  
Gives [Bad, in testing]  
results

# train Accuracy > test accuracy.

## # Underfitting

- train accuracy = 60%  
→ Train Error = 40%.
- test accuracy = 90%  
→ Test Error = 10%.

[Model] [Performs bad in training]

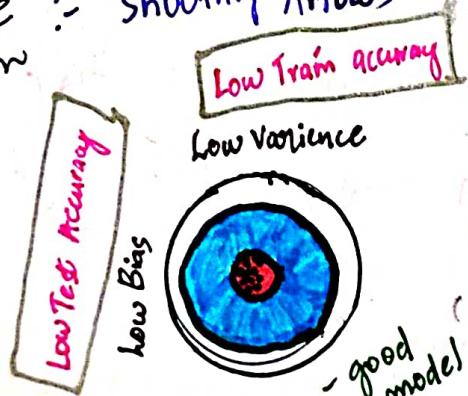
Data. But

performs [Good in Testing data]

# train accuracy < test accuracy.

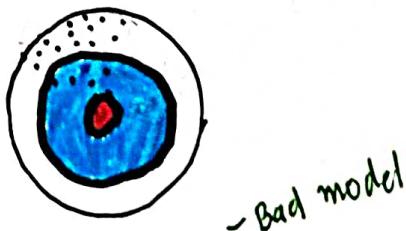
⇒ Bias/Variance Trade-off → difference

Example :- Shooting Arrows



High Variance / High Train accuracy

- Overfitting



Variance = Distance b/w Data Points / Train Error

Bias = close To Target / Test Error

# Overfitting  
~~~~~\*~~~~~

Variance = train  
Bias = test

→ High variance = High train Accuracy

→ Low bias = Low test Accuracy

# Underfitting  
~~~~~\*~~~~~

→ High bias = High test accuracy

→ Low variance = Low train accuracy

# good model  
~~~~~\*~~~~~

→ Low Bias = Low Test accuracy

→ Low Variance = Low Train accuracy

~~1/4/22~~  
1/4/22  
5:00 pm