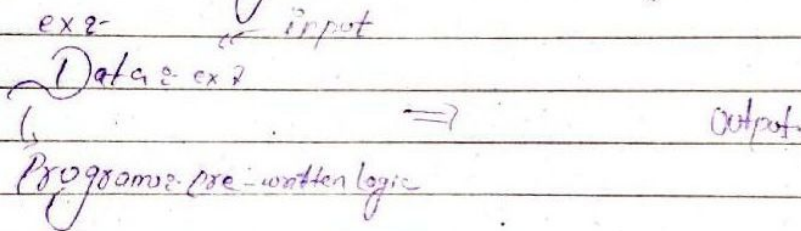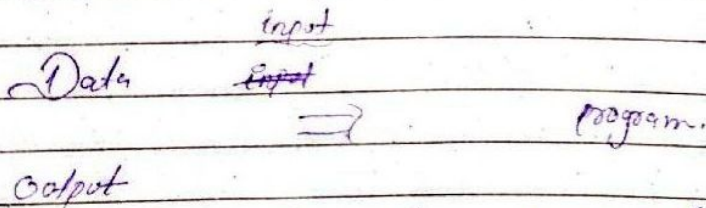=> ML : is the science of getting computers to learn and act like human do, and improve their learning over time in autonomous fashion, by feeding them data and information in the form of observations and real-world interactions.

=> How ML is Different from Traditional Programming?

=> Traditional programming :- it is basically the manual process. in this one human or one programmer built one program. and that we feed to the computer and then it generate the output.

ex :-          input

Data : ex ?

⇒                                    output.
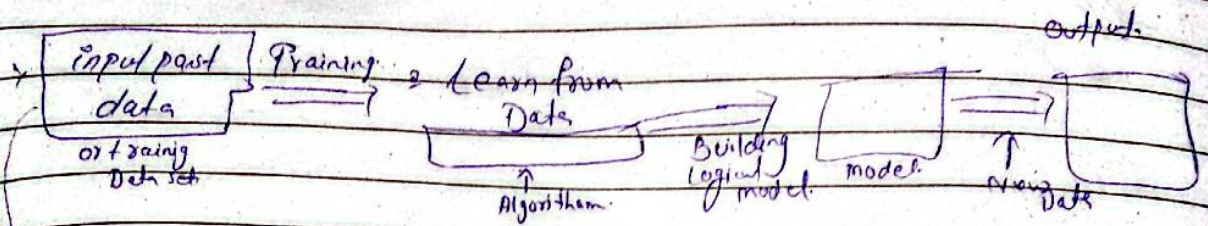
Programme :- pre-written logic

=> Machine learning :- it is basically automated process. so in this we give input as data and output and through the help of Data and output our machine return the program. or logic or model.

                input
Data        input

⇒                          program.

Output

so how machine learn by it self so the simple answer is Data. it is like a fuel.

⇒ How it works.



This input past Data set is feed into the Algorithams. and what is Algorithams? Inside the Algorithams there are set of rules which tells How you this that have to be processed.

2nd phase: In training phase the machine learning Algorithm. in which we feed the data into the Algorithams. that Algorithams learns from that Data's patterns and it does classifications. so in training Phase our Machine learning Algorithim finds the Patterns from that Data, and then it will classified it.

so it is known as training phase from, for machine learn from data

3 phase: Testing phase. so in this we provid New Data sets To the Machines, those data sets which they have never seen before. so what-ever Machine had lean in previous phase and find out patterns and classifications and other things; so it will use all learning in new Data set. And then it will generate you output and provide you the predictions. so here the thing is if the more you have data the the more your model will be learn and then you will get more accurate answer.

→ So this is how Machine learning Model works.

⇒ classification of Machine Learning / Types.
1. Supervised    2. Unsupervised    3. Reinforcement
   Learning       Learning           Learnings

⇒ Data wrangling : Is the process of cleaning, organizing and preparing data for use.

⇒ What is Dataset?
it is the collection of Records.
and what is Record?
A Record is like a single row in a table or dataset. Each row represents a single record, containing specific information or Data points.

ex :-

| ID | Name | Age |
|----|------|-----|
| 1  | John | 30  |
| 2  | Jane | 25  |

in this example, there are
2 records (rows)

Record 1 : ID = 1, Name = John, Age = 30
Record 2 : ID = 2, Name = Jane Age = 25.

→ Note : and that type of Data or Data Sets definitely will be stored in file. and that file can be different file formats.

ex : most commonly use Formats.
1. CSV : comma seprated value.
2. JSON : Javascript object Notation. in this type of format we store the value in the key value formats. like tree like structure.
3. Excell format. and so many other format

⇒ **Types of Datasets:**

1. Numeric Data (Quantitative) → ex Height weight price of house.

2. Categorical data. (Qualitative) → divide into the groups either yes or no T.F. 0.1, 2, 3, 4

3. Ordinal data.    (Note: categorical Data mai aa sakti Numerical

4. it is the combination of both. (Data bachi weerdi ahy. ex T=1 and F=0

ex: if you have order something and then that App wants from so you rating then ex 1 to 5 1 is like poor quality 5 is like excellent quality so in this case, boke thing exist.

⇒ we can get or download data sets from:

1. kaggle Datasets    2. Amazon datasets

3. UCI Machine Learning Repository.

4. Google's Data search Engine.

5. Microsoft Datasets.

⇒ what is Data preprocessing in ML:

↳ Data preprocessing is a process of converting your raw Data into suitable form.

Means the data we get from the any source Most probably it is not clean format. Inside that May be Noise, Missing value occur so by using Data preprocessing we set our data in suitable format.

⇒ Data preprocessing steps involves:-

1. Getting Dataset. 2 importing libraries.

3. Importing Datasets 4 finding Missing values.

5. Encoding categorical Data.

6. spliting Dataset into Training and Test set.

7. Feature scaling.

⇒ step1 create the new file and import your .csv file.

⇒ step2: import libraries.

→ Import numpy as np.

→ Import pandas as pd.

→ numpy library is used for the scientific calculation

→ pandas library is used for to manage the data set.

numpy helps you store and process numbers fast, esspicially when dealing with arrays or tables of Number, like in data science, machine learing or eng.

⇒ Using panda we will display the files like.

step3 dataset= pd.read_csv ('data.csv')

↳ it is just the variable.

dataset # it will display that dataset.

step4: we have to find which is dependent and which is independent.

exe-

$x + y = z$

x+y is dependent.

and z is dependent.

so. we have data set in which we have three independent variables.

country  Age  salary

and the dependent variable is purchased.

so what we have to do is the seprate out them.

→ so country, age and salary we are going to store into the X.

exe.

```
dataset = df.read_csv ('data.csv')
→ X = dataset [ ['country', 'Age', 'salary'] ]  # two dimeshion array
```

X → it will display in tabular form, but if we need array then [[ ]].values

so this is for independent values now similarly for the dependent variable.

```
Y = dataset [ [ "Purchased] ].values.
Y
```

→ so our next step of Data preprocessing is to find out the missing data from the data. like in Age  NaN and salary NaN both have the Missing Values so

so it is very important to handle the Missing values.

There are two Methods.

1 is to remove that particular rows or column from the data set.

but it is not much effiecient, it may arrize data loss issue.

Yeah we can use 1 method when in particular rows or column there Must be a 70 to 75% value Missing, then who can Delet or remove it, or Drop.

→ but in give scenario there is only one value is missing on particular rows and coulms.

→ so

the second way is the to replace the particular value to find the median or mode or mean then we will replace it.

so for this we have to import one library.

→ from sklearn.impute import simpleImputer

Now call the simpleImputer() class.

// simpleImputer(): and the pass the value:

imputer = simpleImputer (missing_values =npnan, strategy = 'mean')

↳ it is just variable.

mean is the value. by default we can change it median or mode

now

imputer.fit ()

so

imputer = imputer.fit (n[:, 1:3] ) #this will take all rows and columns wherever the missing values is occur. this means it will takkens the all rows

Now, other method.

n[:, 1:3]= imputer.transform (:, 1:3 )

↳ all rows ↳ 1 to 3.

↳ what it will do what ever the missing value is it will replace them with the means

→ ✱ x # it will now display and fill NaN will mean value.

=> Next step is categorial data.

it is known as categorial encoding data

categorial data hamari wo data hoti hai jis mai labels hoti hai like country and purchased these are labels not like France no the numeric values.

spanich. Yes

-> so catogerial encoding data kiya hota hai yai jo hamary datasets mai country and purchased hai us ko hum numeric mai convert karygai

os hum is liye karty hai because hamary jo most of the algorithams hai woo categrial data par kam nhi paaty. has

ban khuh hee algoritham hai like decision tree jo is type ko easily handle kar sakta hai

-> so we have country and purchased variable. so for country

-> step1: impart module, named label encoder. jo hamary preprocessing library ka part hai

so:

```
from sklearn.preprocessing impart LabelEncoder.
label_encoder_x = LabelEncoder()
n[:,0]=label_encoder_n.fit_transform (n[:,0])
```

→ all rows of n.
→ country exisist on 0

n // it will show in numeric.

⇒ Now Dummy encoding :-
like

| country | France | Germany | spain |
|---------|--------|---------|-------|
| France  | 1      | 0       | 0     |
| Spain   | 0      | 0       | 1     |
| Germany | 0      | 1       | 0     |
| Spain   | 0      | 0       | 1     |
| Germany | 0      | 1       | 0     |

→ so dummy encoding mai fast 0 and 1 values hoti hou.

so for that

from sklearn.preprocessing import OneHot Encoder.
onehotencoder = OneHotEncoder()
n = onehotencoder.fit_transform (Dataset.country, values.reshape(-1,1)).toarray()
n   // it will display
// Now same for Y
labelencoder_y = Label Encoder ()
y = labelencoder_y.fit_transform (y.)
y:_____

→ Now train and test:
from sklearn.model_selection import train_test_split
n_train, n_test, y_train, y_test = train_test_split(n,y, test_size=0.2,
n_train                                           random_state=0)
n_test.
y_train
y_test.

⇒ Now last is future scaling. ex in our dataset
if there is a data who have huge magnitude
and length then we have to made it Ho the
same scale.