

④ Download R & setup google collab.  
ctrl+enter  $\Rightarrow$  to run a selected code

Vario

3 window in R studio

(A) Environment

(B) Files

I will use VS code

(C) Code editor

# Here we start with machine learning

④ machine learning processes

- (A) Data Preprocessing
- (B) Modeling
- (C) Evaluation

this hierarchy is called  
as machine learning  
process.

these three are main core of entire machine learning and  
we can do anything in ML on basis of these three steps

These above stages can further resolved in another sets

# Data Preprocessing

- (I) Import the data
- (II) Clean data
- (III) Split data in training  
and test sets

# Evaluation

- (I) calculate performance metrics
- (II) Make a verdict  
about model

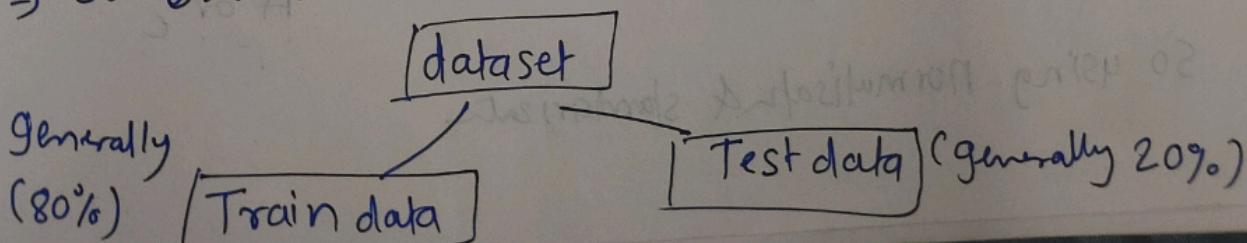
→ Evaluation helps us to find out about  
performance of model

so this is basically  
a structure each  
machine learning  
developer must gone  
through

# Splitting of data into train ~~test~~ & test data

splitting of data is important in order to find out information  
about the model

$\Rightarrow$  So overall dataset is chopped in two pieces



We use 80% Train dataset to build a model & rest 20% is used to test the trained model.

⇒ and such like that we can find out predicted values for test set and check it with actual result and by this we will be able to evaluate model.

### #Feature scaling ] (scaling is good & important)

Feature scaling is always apply to columns of dataset and you will never apply feature scaling to row.

Feature scaling can be done with:

① Normalization

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Above formula will be apply to each value in column & as a result you will have value betn [0:1]

Range of values : [0:1]

Q) 3 persons

		Salary	Age
1	A	70000	45
	B	60000	44
	C	520000	40

Find out with feature scaling whom with Person B is like more likely A or C

So using normalization & standardization

Person	Normalised salary	Standard salary	Normalised age	Standard age
A	1	0.731	1	0.925
B	0.44	-0.052	0.8	0.463
C	0	-0.68	0	-1.4

$\mu = \text{mean}$

$\sigma = \text{std. deviation}$

$$\bar{x}_{\text{salary}} = \frac{182000}{3} = 60,667. \quad \bar{x}_{\text{age}} = \underline{\underline{43}}$$

$$\text{Variance}_{\text{salary}} = \frac{\sum (x - \bar{x})^2}{N} = \frac{(70000 - 60,667)^2}{3}$$

$$\sigma_{\text{salary}} = 12754$$

$$\sigma_{\text{age}} = 2.16$$

$$\sigma^2 = \frac{(70000 - 60,667)^2 + (60000 - 60,667)^2 + (52000 - 60,667)^2}{3}$$

$$\sigma^2 = 162666667 \quad \underline{\underline{\sigma = 12754}}$$

standard.

standard.

standard.

→ we don't want to find out difference on basis of subtraction  
only we want more optimised technique in our case  
it is normalisation & standardisation

⇒ In each column value share same properties

⇒ In row there can be different value for different properties.

By these values we can find out with which value the similar

↳ so following frogmi dilution man

so dilution will be done and we will add all the

## # Preprocessing In Python

[Python]

as main stages in ML

- (A) Data preprocessing
- (B) Modeling
- (C) Evaluation & optimisation

So we can do it all in R or in python but here we are going to do it with python.

→ data is not uniform also not normalised  
So it contains regularities like empty space or wrong value  
In order to remove them from system we use ML.

# Independent variable : what you have

# dependent variable : Variable which we want to predict.

### # Tasks we have to do in preprocessing

- (A) Taking care of missing data
- (B) Removing irregularities & redundancy
- (C) Encoding data (mean verbal data to numbers)
- (D) Splitting data in Training and Test
- (E) Feature scaling (normalising) Standardizing data

[Imp] we divide the dataset in preprocessing state itself.

→ If you preprocess data in very well manner it becomes easy to improve efficiency of model

{# Start with data preprocessing}

### (A) Importing libraries / packages

use keywords like `import, from, as`.

Ex Import pandas as pd

from matplotlib import pyplot as plt

⇒ such like that we can do to use the functions inside the package which can help us in ML & further steps

## ③ importing dataset

for importing file should be at appropriate position of file system  
we mostly use pandas to import data set it has  
functions like read-csv(), read-excel, read-gbq & the  
type of file you want

Ex: `data = pd.read_csv('file name with extension').`

`Pd.read-csv()` give return in a dataframe format  
you can use `Print type(data)` to check datatype of data

**#Important**

You also have to create   
A Matrix of features  
B dependent variable vector

① Vectors are row

② Features are columns other than last

The final data should be on

⇒ last column provide info about verdict mean final  
result which is dependent variable (which we have to predict)  
→ All other feature other than final col are independent

so we have to create 2 entities

A dataset without last column (matrix of feature)

B only last column (dependent variable vector)

In ML on basis of matrix of feature we have to predict  
dependent variable value

Ex: `df.columns["name"]` ⇒ this will give all data for  
a col (feature "name")

so for matrix of feature to

~~ref~~  
you can use iloc function where you can specify  
rows or column which you want

`df.iloc[rowstart:rowend, colstart:colend]`.

so with iloc it becomes easy

so

we want all rows for  
① matrix of features → all col except last  
 $mof = \text{data}.iloc[:, :-1]$ ; → this will work regardless  
of all rows size of data

~~duv =~~ ② for dependent variable vector  
 $duv = \text{data}.iloc[:, -1:]$ ;  
↑ only last col,

now  $mof = mof.values$  &  $duv = \text{data}.duv.values$

this will convert  $\begin{bmatrix} mof \\ duv \end{bmatrix}$  datatype to the ~~array list~~  
humpyarray

using values is not necessary as we can use dataframe  
or it also dataframe can easily converted to ~~list~~

by  $\boxed{\text{data} = \text{list}(\text{df})}$  → numpyarray  
 $\boxed{\text{data} = \text{np.array}(\text{df})}$

values also removes a title of each column like blah blah  
→ which help to focus on data

# mostly the matrix of features and dependent variable  
vector directly used so

$\boxed{X = mof}$  &  $\boxed{y = duv}$  and with these  
you can create your model.

### (C) Taking care of missing data

missing data is a Redundancy in dataset so it is important  
to remove it

$\text{NaN} \Rightarrow \text{Null value}$

NaN is shown in python  
when there is a missing  
value in dataset

⇒ Removing redundancy is

Important in order to make  
model work fine

## ① How to handle redundant data

### i] deleting redundant data.

this is good for large data, but it is not a great as if you have a large data with large redundantness so that is a problem.

### ii] Actually Replacing missing value with average

this is good idea that to fill a empty entry just ~~remove~~ add that entry with a value of average of all from that col (not row obviously)

so how to do it.

for this we are going to use scikit Imputer  
~~so how~~ which rectify nan values & refill it with according to info provided

code like: ~~from scikit import Imputer~~

~~from scikit.impute import simpleImputer~~

~~as Skylearn~~

↑ It is a class

so create instance (object)

Imputer = simpleImputer(); we don't have to use  
new keyword here like we use in java

Simple Imputer take an

so

```
Imputer = simpleImputer(missing_values=np.nan,  
                         strategy = 'mean')
```

This mean we have to enter which type of missing value you want to rectify & which strategy you want to use

and that will remove that value and provide new set

so you created object so now apply it on dataset, matrix of feature and the dependent variable vector.

So now we have connected dataframe & for that  
we have to use fit & transform.

So fit can only apply to numerical columns so you  
have to specify range of dataset with col-

Imputer.fit( $\underline{x}[:, ;A:B]$ ); Imputer.transform( $x[:, A:B]$ )

this is depend on data  
fit will connect imputer  
object with data and  
transform will make  
changes

→ or you can write function for it.

→ same parameters should be entered in fit + transform

the fit transform method returns a value & which will  
be new,  $x$ .

$\underline{x[:, , A:B]} = \text{imputer.transform}(x[:, ;A:B])$ .

↑ we only want change some specific area  
so thus you can remove ~~the~~ missing value

(whole code)

# from sklearn.Impute import SimpleImputer  
# Create object

Imputer = SimpleImputer(~~missing\_value = np.nan,~~  
strategy = "mean")

# now sticking object with dataset & modify it

Imputer.fit( $x[\text{rowrange}, \text{colrange}]$ )

↓ same value to be inserted

$x[\text{rowrange}, \text{colrange}]$  = Imputer.transform( $x[\text{rowrange}, \text{colrange}]$ )

↓

So, thus how you can remove missing values

the fit will get values to  
modify and then  
use transform to make  
it do

fit - transform  
Pair (P)

Here you  
have to  
write  
this  
code  
appropriately

### ③ #Encode categorical data

MLM → machine learning model

So If some column has some set of repeated entries, it become difficult for MLM So we can encode such data in some encoding Ex our data

Country	Age	Salary	Purchased
France	44	72000	No
Spain	27	48000	Yes
Germany	30	54000	No
Spain	38	61000	No
Germany	40	63778.	yes
France	35	58K	yes
Spain	38.77	52K	No
France	48	79K	yes
Germany	50	83K	No
France	37	67K	yes

dependable

In purchased  
if of yes  
repeating

So in country

France
Spain
Germany

repeating multiple  
time

Independent

So while encoding categorical data it can be a

① dependent variable

② Independent variable

And for both we have different way to approach.

① For independent

so approach can be 0: France, 1: Spain, 2: Germany

But this can lead model to some discontinuities so in order to avoid it all we can do is to create

We are going to use one hot encoding and it is very helpful in categorical variable

like for drink

such like that.

- xtra small, small, medium, large, extra large

# the first dependent can be binary output but first we learn for independent  
[Scikit-learn  $\Rightarrow$  sklearn]

Onehot encoding from scikit-learn

#

```
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder
# create object both class
ct = ColumnTransformer([Transformer, remainder])
```

```
ct = ColumnTransformer([Transformer=[('encoder', OneHotEncoder)]])
```

```
ct = ColumnTransformer([transformer=[('encoder', OneHotEncoder(), [0]), remainder="passthrough"]])
```

What should we do for remaining set

Braces

Which column you want to apply transformation

no. of different categories = new

no. of different cols

3 cols as only have Spain, Germany, Italy

# now use that object and connect to it

the columntransformer has fit\_transform which do both fit & transform at once.

X = ct.fit\_transform(X) → do change to.

# this is not numpy array or

so X = np.array(X)

X = np.array(ct.

fit\_transform(X));

And thus how you can convert a categorical encoding

It divide create multiple column so If con is game vector will be  $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$  & such three column

no apply for categorical transformation with  
dependent variable

# we can transform dependent into 1 & 0 and  
for that we have label encoder so for a string  
it convert into number

means:  $\begin{cases} \text{Yes} = 1 \\ \text{no} = 0 \end{cases}$

① we have `sklearn.preprocessing` import LabelEncoder  
and `sklearn.compose` contains transforming.  
Code do same like create tuple and use ColumnTransf

$\Rightarrow$  `LabelEncoder` do not need

`from sklearn.compose import ColumnTransformer as CTR`  
`from sklearn.preprocessing import LabelEncoder`

`LE = LabelEncoder()`  
and `le.fit_transform(y);`

Let column has  
fit-transform

So as we splitted X, y thus how it make us ease

so use scikit for LabelEncoding, OneHotEncoding  
ColumnTransformer and also Impute transform

Fit & transform different

MOF

X  $\Rightarrow$  should be in numpy array

y  $\Rightarrow$  should be  $\Rightarrow$  should not be no much issue

#### ④ Splitting data in Training and testing dataset

so firstly feature scaling is very helpful for normalisation & standardization of data so there is debate that when to use feature scaling before splitting dataset or after splitting.

#### ⑤ Splitting first and then feature scaling.

Feature scaling is used to scale data in some range in order to avoid one specific data to overoccupy all ~~this~~ logic in model

(IMP)

⇒ we want to keep Test Set unknown to training set so in order to hide attributes of Test set like mean, variance, std. deviation we have to first do splitting and then go for feature scaling differently for train & test set

#### ⑥ Basically to avoid info leakage from training set

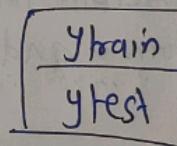
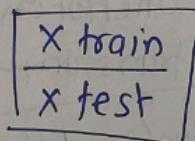
#### ⑦ Splitting data

We are going to use scikit learn in order to split data

so we are going to make 4 sets.

2 for matrix of feature + 2 for dependent variable vector (y)

That are



⇒ And all of them are used in building future machine learning model

so

from sklearn.model\_selection import train-test-split.

So we get 4 outputs

X-train, X-test, y-train, y-test

matrix of feature

dependent variable matrix

= train-test-split(  
Parameters → X,y,  
test\_size → Size = 0.2, random\_state = 1 )

so

Train-test-split (take matrix, take dependent, of feature 'variable vector', test-size = 0.2, random-state = 1)

$$0.2 = 20\%$$

$\Rightarrow$  so just use the test-train-split & you will get four dataset as output.

$\Rightarrow$  If you not provide test-size it will split in 70% Train 30% Test

$\Rightarrow$  It should be in mind while dealing with it.

$x_{train}$  and  $y_{train}$  are related and  $x_{test}$ ,  $x_{tgt}$  and  $y_{tgt}$  are related.

### # Feature scaling

Feature scaling used to avoid overweighting some values over another

So it is not mainly used in all cases but very useful in some of them

# Approach  $\Rightarrow$  good most of the time

Ⓐ standardisation

$$x_{stand} = \frac{x - \text{mean}(x)}{\text{standard deviation}(x)}$$

Range  $\Rightarrow$  -3 to 3

↳ good for specific

Ⓑ Normalisation

$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

Range  $\Rightarrow$  0 to 1

$\Rightarrow$  If your data has normal distribution normalisation will be very good thing to do

$\Rightarrow$  standardization work good in all time

## #Code scikit

We use standard scalar class

```
from sklearn.preprocessing import StandardScaler
```

```
sc = StandardScaler()  
↑ do not need parameter
```

dummy values like values in 0 we don't have to apply standardization here as value are already betn -3 to 3

0	0	1
1	0	0
0	1	0

not good to undergo standardization

so we have to exclude it from data

⇒ only apply feature scaling for numerical value which are high in value like in thousand and avoid to apply it for value in  $\Theta -3 \text{ to } 3$

so

# whom to apply

only apply standardization to ~~test and~~  $x_{\text{test}}$  and

$x_{\text{train}}$

so

$x_{\text{train}}[\text{range}] = \text{sc}.fit\_transform(x_{\text{train}}[\text{range}])$

~~$x_{\text{test}}$~~

and

for  $x_{\text{test}}$

$x_{\text{test}} = \text{sc}.transform(x_{\text{test}}[\text{range}])$

[range]

you can't use fit transform so just use same scalar to transform ↑ don't fit to  $x_{\text{test}}$  will

sc that will lead to information leak