

P^{int}(input("Enter a string:"))

Day:

print(p)

1. 10
D - 9

Date: / /

① format sa values fix kr sktey hai

print(f"Value={1}, Value={2}"). format(3.1, 2)
"Value = 3.1, Value1 = 2"

② ⇒ List ek data type h.

Q1st = [2, 3, 7, 9]

print(Q1st[1:4]) "print(379)"

⇒ nested List

L = [2, 3, [1, 2, 4], 5, 6]

⇒ print(L[range(1, 10)]) "1-9"

③ Sorting

We cannot sort string directly. First, we have to convert them into list than we can sort them.

S = 'Tarkhi is good'

S2 = list(S)

A/HY

S2.sort()

print(S2)

Ais ma char ma bara
bara.

S2.sort()

S3 = ''.join(S2)

print(S3)

Ais ma wapis
join kr di
string m arriver
karay gi.

④ L.append("A") "It will add
to the last of the
list."

L.insert(1, 'B')

To handle
parameter
ma & ha
Value ka gt.

↳ remove(A) You want to like A

→ Ais key andy bs wohi aye False
ga jo remove hua hui, koi index
wah scene ma hui.

length calculate किसी का लंबाई

Day: Print(`len(strme jis ki calculate karne).`) Date: / /

Agi index ka pta. ky tu remove.

def & [2]

→ tu second index removed.

⑤ Tuples:-

$$p = 1, 2 \quad \text{or} \quad p = (1, 2)$$

$$x, y = p$$

print(y)

print(p[1])

→ Tuples can not be changed.

⑥ Dictionaries:-

key values.

dict = {
 122 : 'Rabbi',
 127 : 'Tashi',
 2042 : 'Shubu'
}

print(dict) "Full dictionary Print.(dict). me"

print(dict[127]) "Tashi" \Rightarrow Agi mad aye 127 nth.

print(dict.get(2041)) \Rightarrow Agi ma none aye ga gya.
2042 dictionary me nth. hogi

print(dict.keys())

print(dict.values())

print(dict.items())



~~Notes~~

Question:-

(1) colon op + AKA ka
Dost ma extra try.

y, nai.

(2) def fun(x):
return

print [Exxx.2. for u in L.]

⇒ "if else"

⇒ Loop:-

```
for i in [1, 2, 3, 4, 5]:  
    print(i)
```

⇒ Range:-

```
for con in range(5)  
    print(i)
```

⇒ *for dictionary*
for key, value in params.items()
print(key + " = " + str(value))

⇒ Example:-

```
for int i=0; i<num; i++)  
    cout << i << a[i]
```

Print num
and names.

Same in
python

for idx, x in enumerate(*list*)
 print(idx, n)

⇒ *L1 = [n*x.2. for n in range(0, 5)]*

⇒ while i < 10
 print(i)

⇒ A Keyword for functions called def.

⇒ int Square (int n) // C++

⇒ # Python

```
# def fun(n):
    return n**2, n**3, n**4
```

```
v = square(2)
print(square(n))
```

⇒ def function():
 """ Hey """
 help(function) # it tells what is happening
 # in the function

⇒ def Square (n,y, debug)

Varieties

Square(debug=True, y=1, n=2)

if you write name than

can
only do Date: / /

small things like addition,
multiplication?

Ex: square / fint y^2

$\Rightarrow f_1 = \lambda x: n: x \times 2$

Syntax: lambda argument: expression

we use
this, too
don't use
let

def Square(n):

return n * n

list.

$\Rightarrow \text{print}(\text{map}(\lambda x: n: x \times 2, \text{range}(n)))$
map(function, iterable)

\Rightarrow Range ka Square lae aur list bana dae

\Rightarrow Classes of class, es def __init__(self)

class Point:

this
constructor

def __init__(self, x, y):

Member functions
are always
starts with

self, x, y = self.x, self.y // By default self

$\Rightarrow \text{print}(P)$

ostream Wala kaam / print kرنے کا لیک

operator overloading

Point P1 = 50

P1, P2

print P1

Point P2

def __str__(self)

return f"{{self.x}, {self.y}}"

⇒ Modules :-

import types

dir.

import math

OR

→ from math import *

→ from math import pi

Q.)

→ math.pi

→ math.log

⇒ File :-

files = open("11.1one.txt", "r")

lines = files.read().splitlines()

print(lines)

Qstn :-

⇒ Exceptions :-

if:

else:

raise Exception("life")

try:

normal code.

except:

or

except Exception as e:

→ Ans ma compiler khud hal error bata de.

print("Exception" + str(e))

Exclusive
False

y: → numpy array → famous data science library Date: / /

⇒ Numpy:-

import time import numpy as np
import timeit
%timeit [i+1 for i in range(10000)]

%timeit

↓

for one line

%timeit ⇒ for program

%timeit np.arange(10000) + 1.

jaisy python ma

range by numpy

ma arraayl.

→ Creating numpy array :-

a = np.array([1, 2, 3, 4, 5, 6])
print(np.array(a))

1D → (6,) ∵ 1D array hai jis ma 6 element
hai.

2D → (3, 3) ∵ 2D array hai jis ma 3 rows
aur 3 columns

→ h = np.arange(10)

k = np.linspace(1, 10, 9) means 10 ko
9 parts ma kro.
 $\frac{10}{9} = 1.1$

1. 1.1 2.1

2.1 3.1 3.2

→ import matplotlib.pyplot as plt.

import numpy as np.

y = np.ones(9) → rows.

y = np.zeros(9) "Created an array with 9 zeros"

y = np.zeros((9, 1)) → It's double so write
in 2 brackets.

plt.plot(k, y, 'o')

plt.ylim(0.5, 1)

Exclusive
Falcon

Day:

Date: / /

$\rightarrow y = np.\text{loadtxt}(\text{"file.txt"}, \text{delimiter}=',', \text{print}(y))$

↓
ignores ,
in the text.

for random !! $\rightarrow np.\text{empty}(3,3)$ // We just have to
make an array
Compiler know how
initialize to degg.

$\rightarrow np.\text{identity}(4)$ // identity matrix,
diagonal all same,

\Rightarrow Indexing:-

$a = np.\text{array}([1, 2, 3, 4, 5])$

$a[0:5] \# \cdot a[0:5]$

$a[0:5] \#$

$a = 1, 2, 3, 4, 5$

$b = a[0:2] ==$

True True False False

\Rightarrow Functions:-

~~$a = np.array([1, 2, 3])$~~ \rightarrow ~~np.int8~~ \rightarrow ~~dtype = these are the bits.~~

$a.\text{ nbytes}$

$a.\text{size}$

$a.\text{type}$ $(3,3)$ // tells size (9).

$a.\text{shape}$ $(3,)$

$a.\text{ndim}$ (1) // It's a 1D array go kitni array hy, 1D ya 2D.

$a.\text{dtype}$ \rightarrow ~~Falco~~

Date: / /

a. argmax() // returns index where the max value is.

a. argsort() // sort the indexes
butze sa koy indexes too

$Z = np.zeros(10, np.int32)$

np.ones(10)

$K = p.reshape(2, 5)$

K.ravel() → multidimensional to
1D ma kro doo.
(Flatten kro doon)

np.sum()

np.

Efficiently

~~numpy~~

a.reshape(rows, columns)

Day:

b = $\begin{pmatrix} [1, 2, 3] \end{pmatrix}$ for array.

Date: / /

$\Rightarrow a = np.array([1, 2, 3, 4, 5], np.int8)$, b = $np.array([1, 2, 3, 4, 5])$
print(a+2) # 3, 4, 5
print(a-2) # -1, 0, 1, 2, 3

we want to compare.

Print(a == b) # True, True

Print(a >= b) #

\Rightarrow Transpose.

\Rightarrow Mathematical :-
 $b.T$
 $b.max()$
 $b.argmax()$
 $np.mean(a)$
 $np.std(a)$ \Rightarrow standard deviation.

Pandas:- numpy had limit. so, we use Pandas.

Data Series.

(1D array)

Data frames

(2D array)

\Rightarrow import Pandas as pd.

df = pd.read_csv("data.csv")
↑
obj

`df = pd.Series(a)`

`df = pd.DataFrame(b)`

↓
2D Array.

①

modules (no errors, bhr sa uttarlo)

Built-in

(Pip se install)

External

(People wrote)

r = read

a = appending

② f = open(..., "r")

w = writing

('r' mode is default) f = open("name.txt", 'r')
 l = f.read() OR l = f.write("task").
 print(l)
 f.close().

Modus

① read(r)

② write(w)

③ append(a)

④ create(n)

if file phly ki
toh error.

// Now, we don't use close.

* (With) → say khud hae
close.

⑤ text(t)

⑥ binary(b)

"in case of append."

f = open("name.txt", 'a')

l = f.write("a")

f = open("name.txt", 'rt')

f = open("name.txt", 'rb')

Exclusive
Falcon

Day:

error handling.

Date: / /

(3) Exceptions:-

① `a = input ("Enter a number:")`

try:

```
for i in range(1, 11):
    print(f"\{a} \times {i} = {int(a) * i}")
```

except Exception as e:

```
print("Invalid")
```

```
print("Hello")
```

(2) try:

if

gives

④ Numpy - efficient space

• khud memory handle kar sakte hain

⇒ Numpy Array vs lists.

(more faster)

we can work on

only one data
types. (For big
programs).

⇒ To check dimensions of array

\downarrow

```
a = np.array([[1, 2, 3, 4], [1, 2, 3, 4]])
```

\downarrow
they should
be same

" while printing array , they print in lists
`print(a.ndim)` " prints the dimension
`print(np.array([1,2,3,4]), ndim=10)`

⇒ creating dataFrame from dictionary

✓ import pandas as pd.

① `data = { names: ["Ali", "Haleemah", "Rani"], ages: [1, 2, 3] }`

② ✓ a) `df = pd.DataFrame(data)` or `.pd.Series()`

b) `df = pd.DataFrame(data, index=[1, 2, 3])` we can not put data b/c series is 1D and it is column representation

⇒ extracting data from index

`df.iloc[2]`

⇒ saving to CSV format

`df.to_csv(" .csv")` CSV Created

⇒ reading the CSV file → parsing first col to index

`df1 = pd.read_csv(" .csv", index_col=1)`

A is Se na, phy
1st col first pay
eye jae gr.

Day:

pd.DataFrame(np.random.randint(0, 5, size=10))

Date: / /

" df.loc [0] [Pages]"

" df = pd.DataFrame(np.random.randint(0, 100, size=(100, 2)))

movies = pd.read_csv("imdb.csv")

" movies.head(10)"

displaying fixed number

" movies.loc [0]"

of 1 row.

movies.

if

" df.loc ['split']"

movies. head(1)

movies.

it will display 1st row

feels about
(row, column)

" df. shape"

" movies.columns" // how much columns.

" movies.rename (columns = {

" Votes": "Vote" })" // name

inplace = True

not ho
gaa ga

Age False

new dataset mai

banta

no ga too

no ga gani

pr pechay

" movies.rename (columns = { "Names": "

ek new file

jo prani

se votes

name maps

or ga

if org viral

to change

the column

names

Day: df.info() Date: / /

df.describe() \rightarrow mean, median, mode, std,

✓ for i in movies:
print(i.lower()) \Rightarrow sb ko
small kar de
ga.

✓ movies.columns = [i.lower() for i in
movies]

✓ movies.isnull()

✓ movies.isnull().sum()

✓ movies.dropna() \rightarrow Null values wahe
rakhe nikal ga.

✓ movies.dropna(axis=1)

→ Aisey column
nikal ga ab
jaise main null values

✓ Extracting data from columns.

re = movies['votes']

✓ m = re.mean()

✓ re.fillna(m, inplace=True)

ab agr ahi salary nahi hoga
to remove karte hoga nahi hoga.
Ahi lia mean value nikal
wo aur aur tha values
daad dae ga.

Exclusive
Falcon

```
import matplotlib as mpf
```

Day:

```
import matplotlib.pyplot as plt Date: / /
```

$n = np.linspace(5, 100)$

~~np.linspace~~
 $n = np.linspace(1, 100, 100)$

~~Fig = plt.figure()~~ → free property stores more.

plt.plot(n, np.sin(n), 'o')

It is the range and it prints a graph with base of n-axis as 5-100

we use coloring) in scatter plot

fig.savefig(" -png")

"display a picture.

from IPython.display import image
image(" -png").

capital O

Simple Line plots.

// *matplotlib inline

plt.style.use('')

// fig = plt.figure()
ax = plt.axes()

// plt.plot(n, n+1, linestyle='solid')
or

plt.plot(n, n+1, 'solid')

solid solid

Day:

Date: / /

" plt.plot(u, np.sin(u))
" plt.ylim (-1, 1)
plt.ylim (-2, 1)

" plt.axis ('tight')

('equal')

⇒ " labeling plots.

plt.plot(u, np.sin(u))
plt.title (" ") } write some
plt.xlabel (" ") } must
plt.ylabel (" ") } label
text

plt.legend().
top right
corner
pt

⇒ " Scatter plot :

(2 ways)

① - plt.plot(x, y, 'o', color = "black")

- yng = np.random A's say 'legend' key satth
our key symbol beh

" plt.legend(numpoints=1, fontsize=13)

" -p for polygon.

Day: Co-variance \Rightarrow C.O. \Rightarrow Do variables have
similarities? Date: / /

(2)

plt.scatter(x, y)

"Heat maps".

"plt.colorbar()



side by side

bar bar

giving colors to

"from sklearn.datasets import load_boston"

⇒ Customizing - Legend

Binning \Rightarrow It's on numerical data

↓
to categorized numerical data.

| | | |
|------|----|---------|
| Ali | 17 | (17-20) |
| Shan | 23 | (23-25) |
| Shah | 21 | (21-22) |

(EDA)

dataset from sklearn.datasets import fetch_openml
boston = fetch_openml("boston", data_id=421)

⇒ from sklearn import datasets

import warnings // ignore warning
warnings.filterwarnings("ignore") // ignore warning

n = boston

n = datasets.load_boston()

n.DESCR // tells details

type(n)

n.keys() // last column is target

boston.target, boston.shape // only rows

boston.data.shape // only rows, columns

`n.feature_names` // target variable means result
nai age gr. (All columns named except first one)

`boston-pd = pd.DataFrame(n.adatas)`
`type(boston-pd)`

`boston-pd.head()` // By default 5 print rows.

`boston-pd.tail()` // last 5.

`n.datas.shape`

// Sklearn ma target nai aesa phr array include kya
target hai.

`boston-pd['MPG'] = n.target`

`n = dataset.load_iris()`

`iris-pd = pd.DataFrame()` Same as above
one,

Day:

IDS (continued)

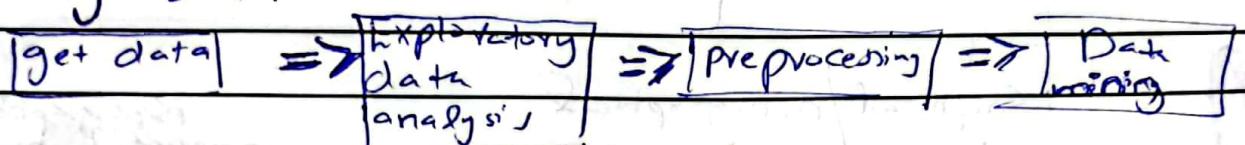
Date: / /

⇒ Exploratory Data Analysis (EDA) (mostly graphical)

Def:- ⇒ Data exploration technique to understand data.
⇒ Approach of data analysis that employs a variety of techniques.

~~Exploratory Data Analysis~~

⇒ Why Data Exploration?



motivations :-

⇒ Select right tool for preprocessing, data analysis and data mining.

⇒ Aim of EDA?

- ① Extract important variables
② Develop valid model.

↓
Proposed explanation
on basis of limited evidence

EDA

- No hypothesis at first.
- Generate hypothesis.
- Uses graphical methods.

CDA

- Start with hypothesis
- Test null hypothesis
- Uses statistical models

⇒ Steps of EDA

- i) Generate good research questions
- ii) Data restructuring → instead of using 2 variables ~~etc.~~, make categories of them.
→ dummy variables for categorical variables.
- iii) Based on research question, obtain descriptive statistics.
- iv) handle missing observations
- v) Decide on need of transformation
- vi) Decide on hypothesis based on researched

Exclusive Falcon

⇒ After EDA :-

- (i) Confirmatory Data Analysis
- (ii) Get conclusions and present results nicely.

⇒ Classification of EDA :-

Graphical/non-graphical

- ↓
- Summarize data in diagram
- ↓
- calculation of summary statistics

Data Profiling = Data calculation & Examining.

Non-graphical

- univariate
- multivariate
- (most bivariate)

one data column
How it varies?

more variables

(before performing multivariate)

Phy Univariate
not real too

Ans Data.

Categorical (Categories)
(Qualitative)

Nominal
(Order not important)

ordinal
(order matter)
Ascending or descending

(unordered data)

level of education

Covered

Rating -

Numerical

Discrete

Counted data

Continuous

measured data

Falcon

Day:

For Text 2

Date: / /

Introduction to Natural Language Processing.

- Human-to-computer language
import nltk.

to write big para.

use "" at start. and " one
at the end.

s = nltk. sent_tokenize(text)

for i in s:

print(s)

for i in s:

w = nltk.word_tokenize(i)
print(w)

text stemming and lemmatization (It's more



accurate)

Change

Change

Changing → Chang

changing → change.

Changes

Changes

Stop words:- English me
jazyky (•)

in nlp it is used as dataset.

(Text dataset
→ Corpus)

from nltk.corpus import stopwords

s = stopwords.words("English")

Variables

Categorical
(Bar chart)

Continuous
Histogram

Uni Variant

import seaborn as sns.

=> iris = pd.read_csv("iris.csv")
iris.mean()

iris.median()

pd.DataFrame([iris.mean(), iris.median()], index=[["Mean", "Median"]])

i = iris.mode(axis=1)

"if I want to change the type"

it
runs
on a
full
dataframe

i.astype("category") or ("int8")

i.describe() (number summary)

i.info()

a = i.cov()

i.corr()

sns.heatmap(i.corr(), annot=True, cmap="Set3")

plt.figure(figsize=(5,5))

sns.boxplot(a)

sns.boxplot(a, orient="h")

// i.cov

// i.corr("Pearson")

// i.corr("Spearman")

//

→ sick

color

(palette="Set1")

"

Day:

Date: / /

- a boston-pd skewness

$n = npo.log(boston-pd.skew())$ (log transformation)
npo sns.distplot(m)

$m = npo.sqrt(boston-pd.skew())$ (square root transformation)
sns.distplot(m)

import statistics

$a = statistics.variance(iris-pd)$

iris-pd.varc() is for every row Variance

Intro to Data Science (IDS)

```
" import nltk
nltk.download()
" from nltk.tokenize import Sent_tokenize, word_tokenize
aText = "...."
print (Sent_tokenize(Text))      "Word tokens
                                         v
                                         list of words.

" from nltk.text import Text
Text(aText)           "Ab tokens ke
" t.count (';') my   text me convert.
                                         If you want to count
                                         the words.

" t.vocab()    " text kay sati yaa batre ga
                                         v
                                         kay kitni baar repeat.
according dictionary bani dae ga.

" from nltk.corpus import stopwords      "data set
                                         in nltk
                                         set (stopwords.words('English'))      "corpus.
                                         v
                                         ya unique
                                         words ko
                                         la khaye
                                         ga.

" filtered_words = []
for i in word_tokenize:
    if i not in stopWords:
        filtered_sentence.append(i)

" Stop.words.fileid()      "languages in stopwords.

" Stopwords.raw('arabic').replace ('\n', ' ')
                                         v
                                         it print karey
                                         Arabic language kay stop words.
```

- ① `// tokenize`
- ② `// lemmatization`
- ③ `// stopwords.`

```
// from nltk.stem import PorterStemmer
// tokenize w-l:
```

" Ais say
original word
age gs.
ing, ed sbr
Khatam.

`ps = PorterStemmer()`

`e = ['python', 'pythoned']`

`for w in e:`

`print(ps.stem(w))`

```
// from nltk.stem.WordNetLemmatizer
```

" "

`for ...`

~~`print(ps.WordNetLemmatizer())`~~

" Stemmer 3 types.

`from nltk.stem.snowball`

`import`

SnowBall

Same auth.

" " . lancaster

" hananster Stemmer

" " . porter

" Porter Stemmer

Used porter.

Porter improvement → Snowball

" import string

punctuation string ma ha aya
jatay hain ai kia nltk use
na kro.

`p = string.punctuation`

`p.`

"First Limitation then Stopwords"

"yield" \Rightarrow it just saves the value to any huge

 used in for loops.

```
" from textblob import TextBlob "spelling to check know key
```

$t = " \dots "$

$$t \mapsto T_{\mathcal{E}X+B \circ b}(t)$$

~~tc =~~ t1.correct()

11 Ngyambo

`text = "..."`

import nltk

netk.download('punkt') " Ya knol sentence ko correct.
netk.download('Finnish') punctuation lgia dee ge.

$m = n \cdot l + k : \text{FreqDist}(\text{text})$

1) unigrams / biGrams / triGrams

11 predictions may e.g. long distance

Unigrams = nltk.ngrams(tokens, 1).

$Uf = m \ell + k$. Frey Dilt (Unigrams)

for token in list (string of items) :

11 Time Series (As date mention hti by, strong me)

11

```
df_hb8['Date'] = pd.to_datetime(df_hb8['Date'])
```

df_hbl_index // ya start, stop
our step, kn
bucle sei

`df['Date']`)
it is a
date
in excel
by two string
meta, num or int or
time series.

Real green level. \hookrightarrow $RGB \Rightarrow 3$ images hi aisa ma. (3 dimensions)
 \Rightarrow OpenCV
C A.I. ma se hum 1 image nikal letay hain. Phr usi me kaam karna easy.

img = cv2.imread('path') // read the image.

img = cv2.cvtColor(img, cv2.COLOR_RGB2GRAY) // 1 row

fig, ax = plt.subplots(1, 2, figsize=(10, 10))

ax[0].imshow(img) // original image
no phly show karay ga.

ax[1].imshow(graying) // upper gray image ko
2nd number py dikhae ga.

// RGB ma split karne ka fin. 3 variables needed.
a, b, c = cv2.split(img)

import librosa

ft → sample points (interval kitna)
det. → +

\Rightarrow spectrogram \rightarrow visual representation of audio.

\Rightarrow time = date / sample rate.

\Rightarrow fill_between //

[a, b, c] = plt.magnitude

↓ ↓

We had used this
D/C

•- Edge detection → Search pixels where there is sharp change.
→ Sobel and Prewitt filters.

→ Prewitt

•- Sobel

$$g = \text{prewitt_h}(\text{img})$$

•- Sobel

$$b = \text{Pfilters.Sobel}(\text{img})$$

•- Canny → highest accuracy
feature.Canny(img)

Image Processing :-

import CV2

img = cv2.imread("profile.jpg")

grey_img = cv2.cvtColor(img, COLOR_BGR2GRAY)

fig, ax = plt.subplots(1, 2, figsize=(10, 4))

ax[0].imshow(img)

ax[1].imshow(grey_img)

// Split image in colors

b, g, r = cv2.split(img)

fig, ax = plt.subplots(1, 3, figsize=(10, 5))

ax[0].imshow(b, cmap="Blues")

ax[1].imshow(g, cmap="Greens")

ax[2].imshow(r, cmap="Reds")

ax[0].axis("off")

Kernel size

// Blur images

g-blur = cv2.GaussianBlur(img, (19, 19), 15)

↑
larger
size, more
blurry.

↑
blur in width.

↓
std

m-blur = cv2.medianBlur(img, 9)

↓
blur value

fig, ax = plt.subplots(1, 2, figsize=(10, 4))

ax[0].imshow(g_blur)

ax[1].imshow(m_blur)

Image Rotation

angle = 180

scale = 1 // Scale means zoom

h, w = img.shape[:2]

centre = (w//2, h//2) // center of rotation (It is important)

rot_matrix = cv2.getRotationMatrix2D(centre, angle, scale)

rot_img = cv2.warpAffine(img, rot_matrix, (w, h))

Image Resizing

new width = 500

new height = 500

h, w = img.shape[:2]

resized_img = cv2.resize(img, (newwidth, newheight))

Image flip

flip_img = cv2.flip(img, 1)

Data imbalance \Rightarrow Red Black] So many na samjhayga time

Performance metrics

① Data preparation \rightarrow feature engineering \rightarrow data modeling \rightarrow Performance measure.

(Features nikalay hain)

② Confusion matrix (Used in classification) Like KK banday ko Cancer machine Yes. array ko cancer mi, No, Rakin agya agar ulte ho sae, Yes ko No ko NO ko Yes. Ab aisi table ko confusion matrix

| | | positive | negative | |
|----------------|----------|------------------|----------|------------------------------------|
| | | Actual (Columns) | | |
| True Positive | Positive | TP 2 | FP 3 | $\Rightarrow 5$ |
| | Negative | FN 1 | TN 3 | $\Rightarrow 4$ |
| False negative | Positive | 11 | 11 | Small cases minimize (covid Spain) |
| | Negative | 3 | 6 | Big cases minimize (cancer) |

(We have to overcome FP and FN according to the situations)

(Agr faili Raking
hwahyto agya
kund galat
ho gya)

(e.g. FN, galat N kya)

Accuracy

| | |
|----|----|
| TP | FP |
| FN | TN |

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}}$$

Precision

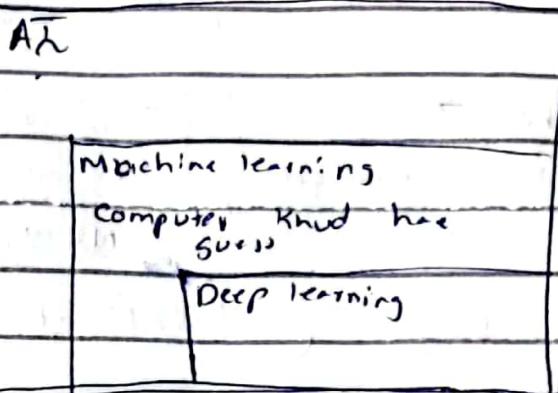
$$\frac{\text{TP}}{\text{TP} + \text{FP}}$$

(Accuracy achieved,
data imbalancing)

Session 8-2.

① Machine learning

"Field of study that gives computer the ability to learn without being explicitly programmed"

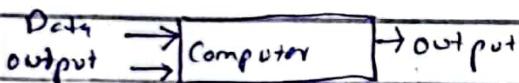


- Start with data $\rightarrow E$ (experience)
- Perform some sort of task $\rightarrow T$ (task) (Task is called supervised learning, classification)
- Validation test \rightarrow Performance measure P or inductive learning

Traditional



ML



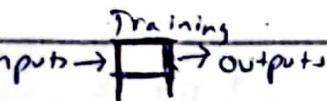
less complex than us
b/c output is known

ML

- detection of anomalies
- association
- clustering

Supervised (Classification)

- learn explicitly
- Output clear
- Predicts future



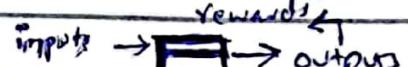
- Medical Diagnosis
- Labeled data
- 2 steps - learn (training)
-- testing
- Decision trees
- KNN

Unsupervised

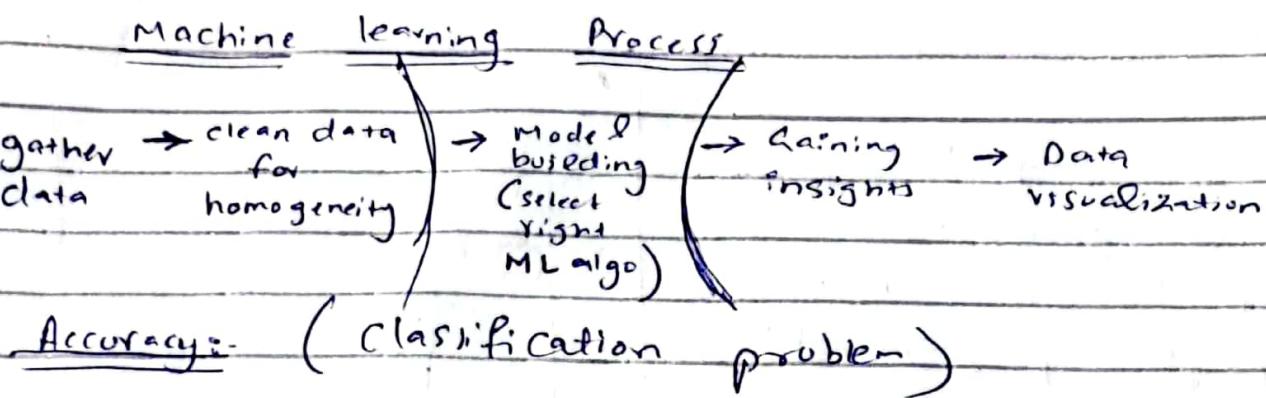
- machine learn
- do not predict



- Fast Students new
- unlabeled data
- K-means clustering



- Reward based learning
- Netflix, youtube
- Recommendations.



Accuracy = No. of correct classification

Total no of test cases

Classification:- (2 things)

- We predict y labels (classes) for input x .
- Yes or No (Fraud or not fraud).
- $y = f(x) \rightarrow$ input
↓
output classification function

Regression:- numerical data (1 thing)

Supervised

- Feedback
- Used for prediction
- Known number of classes
- less complex than \rightarrow

Examples

- Text recognition. (^{training}_{no rahi})
- Spam detection
- Face detection

Unsupervised.

- no feedback
- used for analysis
- unknown number of classes

Examples

Noise removal from datasets!

Decision trees (for classification & regression)

Classification:-

- 2 steps
 - learning step
 - prediction step.



Decision tree algorithms:-

- Topmost node \rightarrow root
- Internal node \rightarrow feature
- Branch ~~node~~ \rightarrow decision rule.
- Outcome \rightarrow leaf
- It learns to partition on attribute value
- It partitions the tree in recursive manner called recursive partitioning.

Entropy

-- measures impurity of input set

$$= \frac{C(P)}{C(T)} \times \log_2 \frac{C(P)}{C(T)} = \frac{C(N)}{C(T)} \times \log_2 \frac{C(N)}{C(T)}$$

$$= \frac{9}{14} \times \log_2 \frac{9}{14} = \frac{5}{14} \times \log_2 \frac{5}{14}$$

Gain:-

$$\text{Entropy (Info)} = \frac{C_{\text{Sunny}} \times \text{Entropy}_{\text{Sunny}}}{C_{\text{Total}}} + \frac{C_{\text{Cloudy}} \times \text{Entropy}_{\text{Cloudy}}}{C_{\text{Total}}} + \frac{C_{\text{Rainy}} \times \text{Entropy}_{\text{Rainy}}}{C_{\text{Total}}}$$

(PFA)

- Attribute selection measures \downarrow known as (ASM).

Splitting Rules.

- Different Metrics :-

Confusion Matrix

Accuracy

Precision

AUC (Area under the curve)

MAPE (Mean absolute % error)

MAE (Mean absolute error)

MSE (Mean squared error)

Actual

| Predicted | | TP | FP |
|-----------|----|----|----|
| Actual | TP | FN | TN |
| | FP | | |

- Accuracy :-

$$\frac{TP + TN}{Total}$$

Total.

TP
TN

- Precision :-

$$\frac{TP}{TP + FP}$$

TP
FP

(How many did we catch)

- Recall or Sensitivity :-

$$\frac{TP}{TP + FN}$$

TP
FN

(How many we caught)

If every person has cancer, recall is 100%.

- If focus on minimize FN, Recall $\rightarrow 100\%$ without precision too bad

- If on " FP, Precision $\rightarrow 100\%$.

-- Specificity :- (Opposite to recall)

$$\frac{TN}{TN+FP}$$

$$\frac{TN}{TN+FP}$$

-- F1 Score :-

Single score that represent Precision and recall.

F1 Score = Harmonic Mean (Precision, recall)

$$F_1 \text{ Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

-- AUC :- (Scale - invariant) (0 - 1 range) (If 100% correct
AUC 1.0
else 0.0)

a) True-Positive rate (TPR)

$$TPR = \frac{TP}{TP+FN}$$

b) False-Positive rate (FPR)

$$FPR = \frac{FP}{FP+TN}$$

-- MAE :-

$$= \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i) \quad (\text{Avg diff b/w original and predicted values})$$

-- MSE :-

$$= \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2$$

-- MAPE :-

$$= \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right|$$

X = Independent / predictor / explanatory

Y = Dependent / response variable. Date 1/1

Linear Regression (for prediction or forecasting)

- Dependence of one Variable on one or more variables (called independent variables).

\Rightarrow Primary used to input to a system

- Dependent & independent variables. Key relationship ka estimate.

- Effect of each explanatory variables on dependent variables.

- Predict value of dependent variable for given ind variables.

\Rightarrow o Dependent are those which are changed due to indep.

- Least Square Linear regression \Rightarrow method for determining Y on basis of X .

Assumptions of Linear Regression Model.

- Linear Functional form
- Fixed independent Variables
- Equality of variance of the errors.
- No multicollinearity
- No outlier distortion
- No autocorrelation of the errors.

Linear Regression Model.

First order linear model

$$Y = b_0 + b_1 X + \epsilon$$

b_0 = Y-intercept

Slope.

b_1 = Slope of the Line

\uparrow Y-intercept.

$$y = mx + b$$

ϵ = error variable.

- Independent change, 1 dependent change
- 2 Independent change, 1 dependent change.

[↑]
predict

Date / /

- If there is only one driver variable, X , it is simple linear regression.
- Model involves multiple driver variables called multiple linear regression.
- If Relationship b/w X and Y is curvilinear, the regression line will be a curved line.
- Greater strength of relationship b/w X and Y better is prediction.

⇒ Least Squares Estimation of b_0 , b_1

Least square estimates of slope co-efficient b_1 of true regression line

$$\textcircled{1} \quad \hat{y} = \hat{b}_0 + \hat{b}_1 x$$

$$\textcircled{2} \quad SSE = \sum_{i=1}^n (y_i - (\hat{b}_0 + \hat{b}_1 x_i))^2$$

$$\textcircled{3} \quad \hat{b}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

\hat{b}_0 = Mean response when x = mean response when x increased by 1 unit.

$$b_0 = \bar{y} - b_1 \bar{x}$$

b_0, b_1 are unknown parameters.

- Regression generates least squares regression line.
- Squaring difference and adding up squared differences across all predictions is called residual/error sum/squares/SSerror.
- Regression generates formula such that SSerror is as small as it can possibly be.
- Minimizing this number minimizes average error.

LUCKY
EXCLUSIVE

ENJOY YOUR WRITING WITH LUCKY PAPER PRODUCTS

Regression in Python

- ① from sklearn.linear_model import LinearRegression
 - ② Do transformation. Use reshape(), set u, y
 - ③ • model = LinearRegression().fit(u, y)
 - fit.intercept \Rightarrow Boolean, if true decides to calculate intercept b_0 .
 (By default true) • if false consider it equal to zero.
 - normalize \Rightarrow Boolean, if true decides to normalize
 (By default the input variables false)
- ↓
- It doesn't normalize input variables.
- model.intercept_ (b_0)
 model.coef_ (b_1)
- ④ Check results of model fitting
 - Obtain coefficient of determination, R^2 with .score()
 $q_1 = \text{model}.score(u, y)$
 - pre = model.predict(u)

Applications

Economic growth

Product sale

Housing Sales

Score prediction.

If $K = \frac{\text{no. of obj}}{1}$, then distance = 0 (case of overfitting)

K-mean Clustering

Date / /

| | n | y | |
|---|-----|-----|----------------------|
| 1 | 1.0 | 1.0 | \rightarrow Mean 1 |
| 2 | 1.5 | 2.0 | |
| 3 | 3.0 | 4.0 | |
| 4 | 5.0 | 7.0 | \rightarrow Mean 2 |
| 5 | 3.5 | 5.0 | |
| 6 | 4.5 | 5.0 | |
| 7 | 3.5 | 4.5 | |

Advantages:-

- Easy to represent
- Can work in multiple dimension

Disadvantages:-

- Time consuming to find optimal number of clusters.
- Time consuming feature engineering

Iteration 1:-

\Rightarrow For point 1:- $n_1 = 1.0$, $y_1 = 1.0$

$$D_1 = \sqrt{(n_2 - n_1)^2 + (y_2 - y_1)^2}$$

$$= 0$$

$$D_2 = \sqrt{(5 - 1)^2 + (7 - 1)^2}$$

$$= \sqrt{(4)^2 + (6)^2}$$

$$= \sqrt{16 + 36}$$

$$= \sqrt{52} = 6.10$$

\Rightarrow For Point 2:- $n_1 = 1.5$, $y_1 = 2.0$

$$D_1 = \sqrt{(n_2 - n_1)^2 + (y_2 - y_1)^2}$$

$$D_1 = \sqrt{(1 - 1.5)^2 + (-2)^2}$$

$$= \sqrt{(-0.5)^2 + (1)^2}$$

$$= 1.12$$

$$D_2 = \sqrt{(5 - 1.5)^2 + (7 - 2)^2}$$

$$= 6.10$$

| Point | D_1 | D_2 | which cluster to choose. |
|--------------|-------|-------|--------------------------|
| 1 (1.0, 1.0) | 0 | 6.10 | 1 |
| 2 (1.5, 2.0) | 1.12 | 6.10 | 1 |

Date / /

Iteration 2:

For Mean 1 :-

$$\frac{1.0 + 1.5 + 3}{3}, \frac{1.0 + 2.5 + 4.0}{3}$$

(Points in
cluster 2)

$$(1.8, 2.33)$$

\downarrow
 y_2

For Mean 2:-

$$(1.12, 5.38)$$

\downarrow
 y_2