

6/06/24

Artificial Intelligence

function definition

Engineering of making intelligent machines &

programs.

function definition

Machine learning :- Part of AI an Ability to learn without being explicitly programmed.

Deep learning :- learned based on deep neural network.

function definition

Types of AI

function definition

Weak AI

(or)

Narrow AI

function definition

Strong AI

(or)

Broad AI

eg: Robots

eg: • Driver-less Car

• Smart home System

• Alexa, Siri, Google Assistant, etc.

Siri

• Playing Ludo or chess in Computer

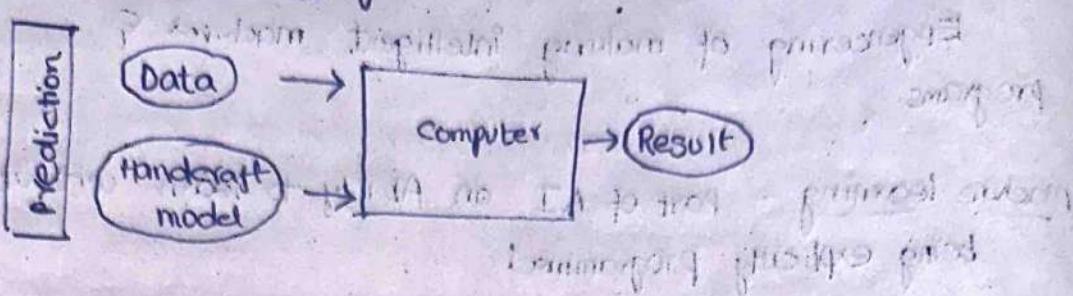
E-commerce Real time Example of AI Zomato, flipkart

* Machine learning:

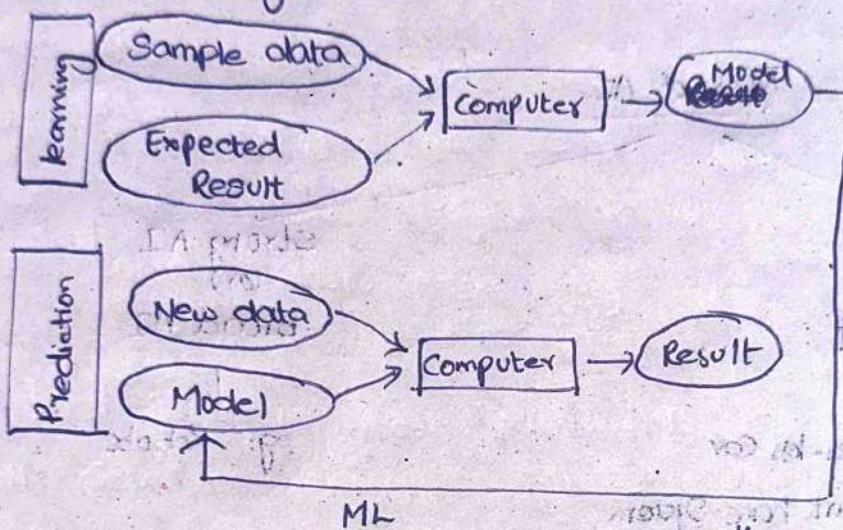
- ML is a way for computers to learn from data & make decisions, without being explicitly programmed for each task.
- Its like teaching a computer to recognize patterns & make sense of information on its own.
- ML uses data to detect various patterns in a given dataset.
- It learns from past data & improve automatically.

Traditional Programming vs Machine learning:

Traditional modeling:



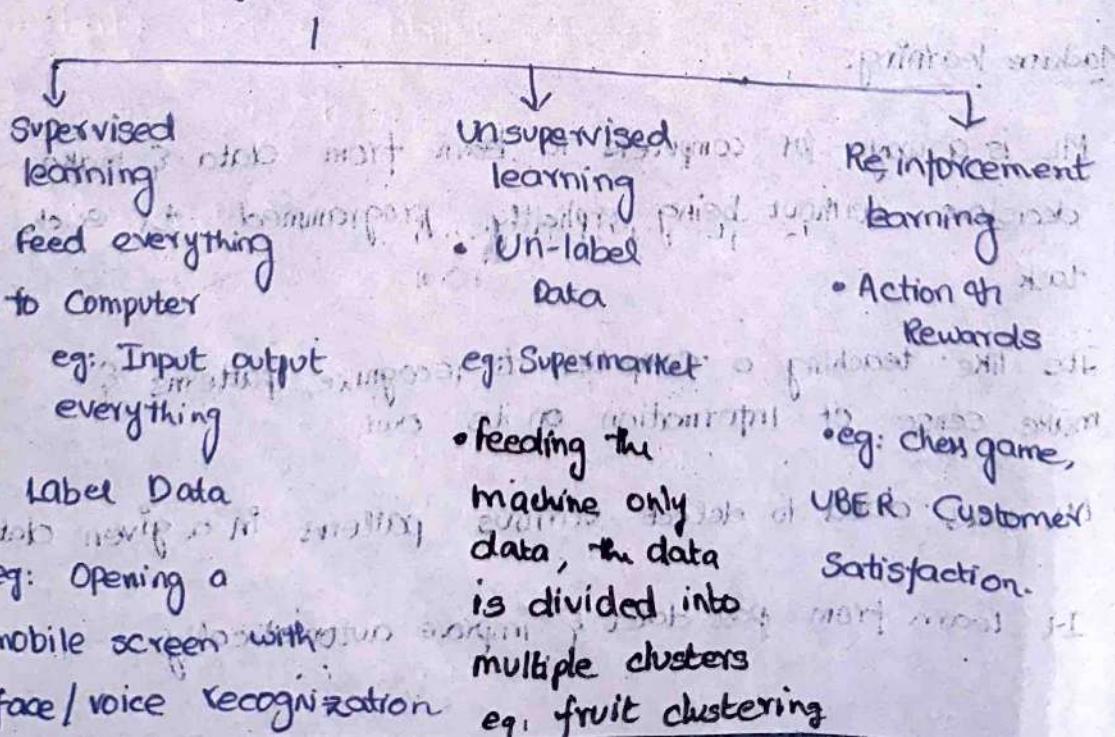
Machine learning



Application to Understand Concept in Real World: Teachable Machine

ML is one of the powerful way to achieve Intelligent
Artificial intelligence

Types of ML



Working Mechanism of ML

- 1) Get Data
 - 2) Clean, Prepare & Manipulate Data
 - 3) Train Model
 - 4) Test Data
 - 5) Improve
- } Data was split into training & testing (80:20)
- (80% = training, 20% = testing)
- difference b/w these two values is error value.

Unstructured data - Consists of different types of rating like eg: chat section of zoom, one may give image reactions etc. one may send answer through text i.e combination of all types of data.

+ Binary Classification:
Representation: $<=50K$, $>50K$ etc. (from adult.csv shared in session)

- To handle large data set we use pandas
- They are simple to use for data structures & manipulating (+) fast & disk

- `head()` - display first 5 entries (default case)
- `tail()` - display last 5 entries (default case)
- Preprocessing - data cleaning process which is an important prerequisite (helpful for the machine to learn data in better manner to produce off data models)
- Inconsistent data set - Some features are not going to contribute much for students employee table.
- data shape - for fetching the details of rows & columns in a table.

select few rows - `(:) .mug .() .onei .ntok`
using precision

Box plot concept:

```
import matplotlib.pyplot as plt #
```

code:

```
import pandas as pd
```

```
data = pd.read_csv(r"path\adult.csv")
```

```
data.shape (48842, 15) (rows & columns)
```

```
data.head(3) (first 3 entries)
```

only first 3 entries

S.No.	Age	Work class	fnlwgt	Education	Ed-N _o	Marital stat	Occupation	Relationship
0	35	Pvt	22648	HS-grad	10	Not married	Fishing	child
1	38	Pvt	89861	HS-grad	10	Married - spouse	Farming	Husband
2	28	Local Govt	336951	Assoc-acdm	12	Married - spouse	Service	Husband

only first 3 entries were printed (data.head(3))
but in case of default case we get first 5 entries.

Now function 'tail' prints last n entries. let's print last 4 entries of data.tail(4)

Similarly last 4 entries were printed.

+ Null values or Missing values

data.isna() - function used to check whether

there are missing values in particular column or

output is binary format

(returns either 0 or 1)

data.isna().sum() - count of those missing values

Employee burnout

Output: data.isna().sum()		Prediction Process
age	0	not known
workclass	0	not known
fnlwgt	0	not known
education	0	No missing
education-no	0	values in given dataset
marital-stat	0	not known
occupation	0	not known
relationship	0	not known

Mean, Median, Mode, arbitrary

- different modes to address missing values
in the given data

* print(data.workclass.value())
 ↓
 oper
ations
 to
be
performed on
 ↓
 only
 worclass
 column
 ↓
 no. of values of
 count of each category

Output:

workclass
 private

33906

self-emp-not-inc

3862

local-gov

3136

otherwise ?

category

state-gov

2799 (21.18%) 312.0

without-pay

[1981] not) 312.9

21

312.0

Not known

in place of? we can replace with other category & use same data.

function:

```
data['workclass'].replace({ '?' : 'NOTlisted' },  
                           inplace = True )
```

```
print(data['workclass'].value_counts())
```

output:

workclass

Private

self-emp-not-inc

33906

3862

local-gov

3136

→ NOTlisted

2799

State-gov

self-emp-inc

1695

Similarly occupation column.

irrelevant features

```
data = data[data['workclass'] != 'Without-pay']
```

```
data = data[data['workclass'] != 'Never-worked']
```

↳ double Array

data.shape

output: (4881, 15)

```
print(data['workclass'].value_counts())
```

output:

workclass

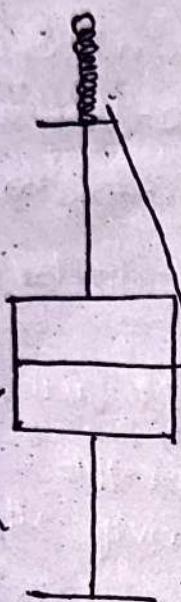
"

table.

Visualization of outlier detection

```
import matplotlib.pyplot as plt  
plt.boxplot(data['age']) # numeric values.  
plt.show()
```

Output:



Box plot concept
only numerical columned values are taken.

* Box plot is only applicable to numeric values

* Bar charts, scatter plots for categorical values.

least information in our data set

These are our outliers in our dataset

to remove outliers in data set

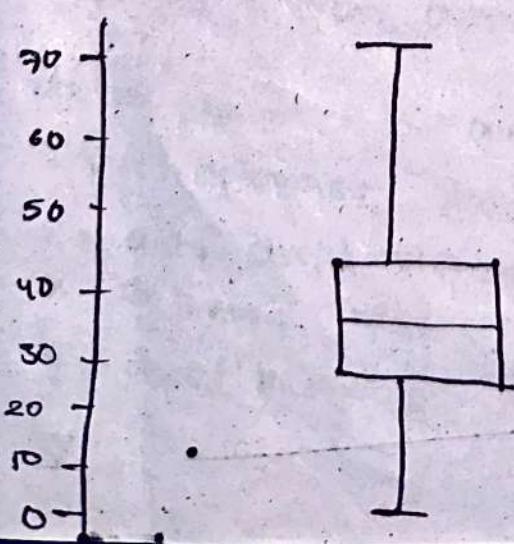
```
data = data[(data['age'] <= 75) & (data['age'] >= 17)]
```

```
plt.boxplot(data['age'])
```

numerical values

```
plt.show()
```

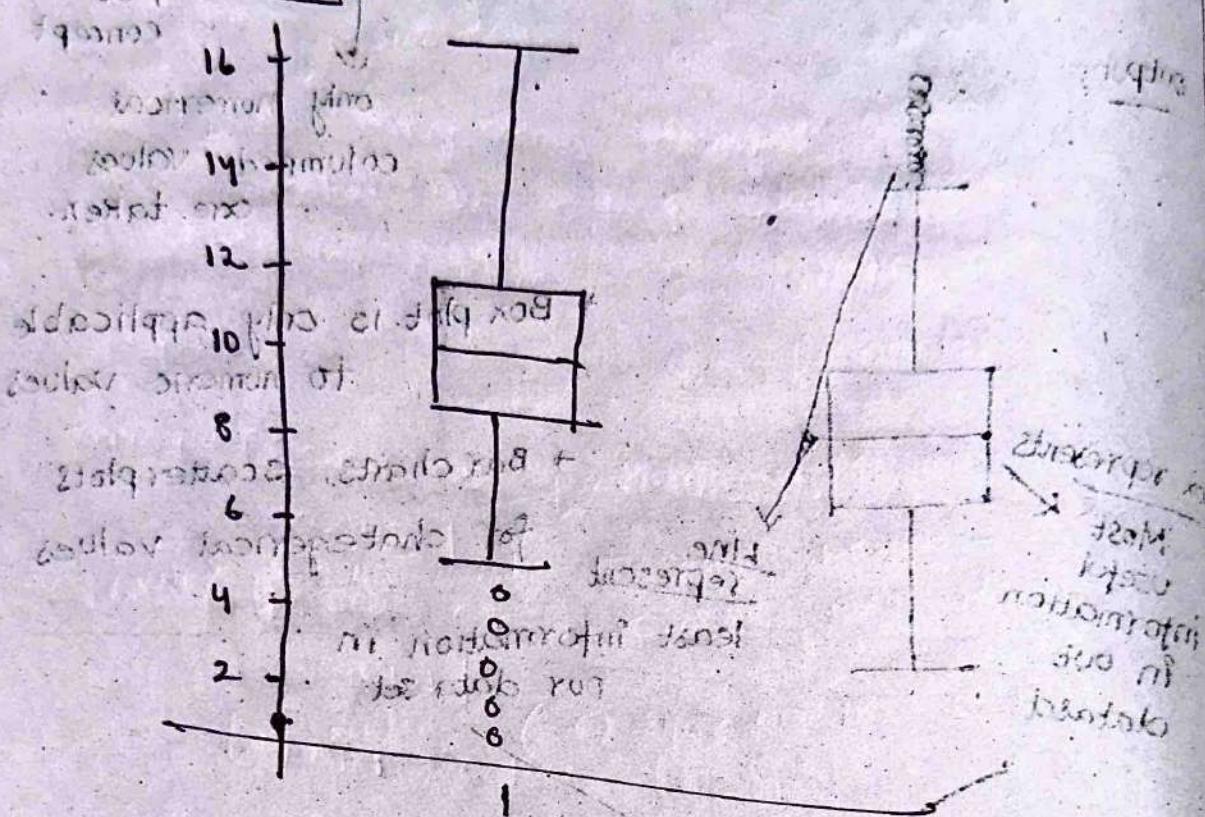
You can observe the outliers are removed from graph



take another column & repeat it

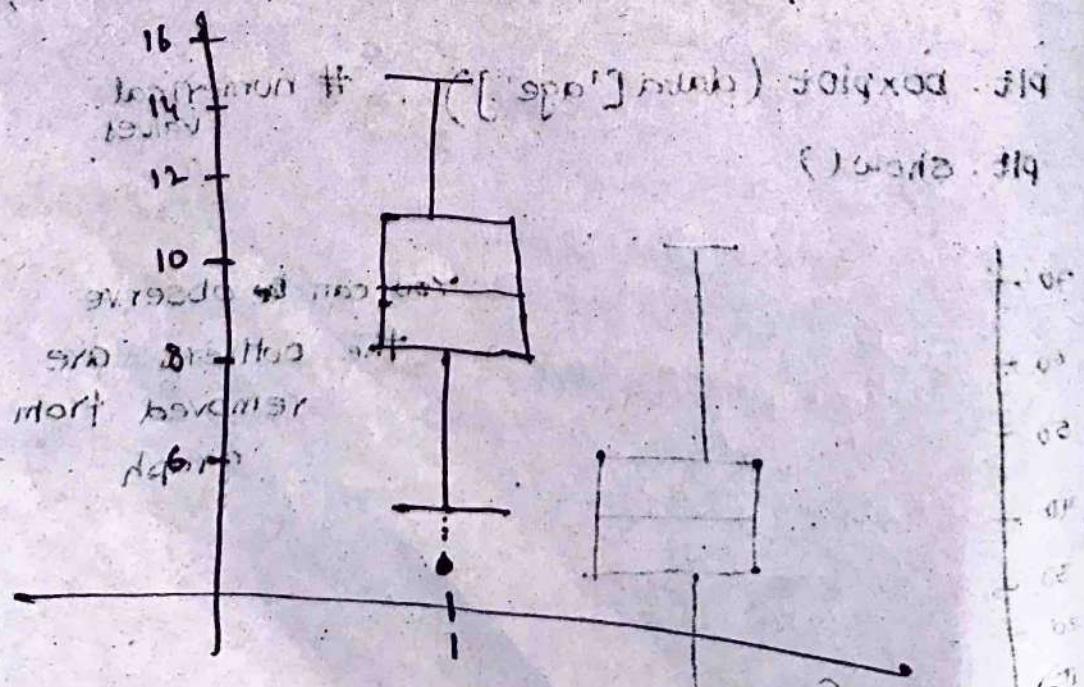
plt. boxplot (data['educational-num']) # numerical values

plt. show.



data = data [(data['educational-num'] <= 16) &
(data['educational-num'] >= 5)]

plt. boxplot (data['educational-num']) # numerical values
plt. show()



* Redundant information

```
data = data.drop(columns=['education'])
```

drop - delete the column from dataset

sparse data - The group with less data
points.

* conversion of categorical values into numerical values.

```
from sklearn.preprocessing import LabelEncoder
```

```
encoder = LabelEncoder()
```

import lib

- Invoking an object

```
data['workclass'] = encoder.fit_transform(data['workclass'])
```

```
data['marital-status'] = encoder.fit_transform(data['marital-st'])
```

```
data['race'] = encoder.fit_transform(data['race'])
```

Similarly other categories

All categorical values transform to numeric values

Now

split input from output for data processing

Output column: last one (here it is income)

Input columns: remaining 14 columns considered

```
x = data.drop(columns=['income']) # input inputs.
```

```
y = data['income'] # output
```

y

* Normalization

```
from sklearn.preprocessing import MinMaxScaler  
scaler = MinMaxScaler() # object  
X = scaler.fit_transform(x) # input  
X # range of 0 to 1
```

All these array values arranged themselves like 0 to 1 to make it easy for machine learning.

Output:

```
array([[0.137, 0.5, 0.1415, 0.397,  
       0.9512],  
      [0.362, 0.5, 0.052, 0, 0.5,  
       0.9512],  
      [0.189, 0.1666, 0.219, 0, 0.397,  
       0.9512],  
      ...  
     ])
```

propagation until final

eg: 1 2 3 4 5 range - 1 to 5
Conversion to Smaller values. 0.2 0.4 0.6 0.8 1 range - 0 to 1

propagation back not doing most things

(initial value) and final values doing propagation in forward pass. X input # (initial) - initial values. X output # X output # [final] - Y

+ training & testing data set: to gain sufficient data

* logistic regression.

```
from sklearn.linear_model import LogisticRegression
```

LR = Logistic Regression ()

```
lr.fit(xtrain, ytrain)
```

ニサツ

logistic Regression → for prediction process
and also accuracy

- Logistic Regression()

Binary

Classification

`predict3 = lr.predict(xtest)`

predicts

→ output can be either positive or negative.

```
array(['<=50k', '>50k', '<=50k', ... , '<=50k'])
```

`dtype = object)`

```
accuracy_score(ytest, predict) # difference of
```

Output: 0.8203125

actual output and
the predicted

↳ 82.1.

Decision tree - hierarchical tree structure based on decision at each hierarchy

- An Machine learning algorithm.
 - At each hierarchy you take decision
 - Binary classification

Employee Burnout Prediction

Linear Regression

is a powerful statistical technique to model the relationship between a dependent variable and one or more independent variables.

What is Employee Burnout prediction?

- Predicting the stress level of a employee
- For an organization it is required to go through the Nature & stress levels of an Employee So it is going to contribute the whole organization's growth.
- So, it is important to address the employees mental health, p stress levels.

What is Supervised learning:

- Basically we will be feeding the machine computer both the inputs & output so it can generate data & apply it to future datasets.
- This model uses Labelled data
- In this project (Employee Burnout Prediction) we will be using Linear Regression Technique.

Dependent and independent Variables

Dependent variable - like an outcome that we want to predict or explain

Independent variable - factors that influence the Dependent variable

- e.g. • Final examination score is DV
• All our struggle, efforts, internal marks, class Assignment, activities that influence our final score is IV

Defining Employee Burnout:

- Employee Burnout is a state of emotional, physical & mental exhaustion caused by prolonged or excessive stress

- 1) Emotional exhaustion
- 2) De-personalization
- 3) Physical symptoms
- 4) Reduced personal accomplishment.

Data Collection & Pre-processing

- online resources : Kaggle

- 1) Data Acquisition
- 2) Data cleaning
- 3) Feature Engineering.

- All these "machine learning models" are based on iterations / data.
- ML is impossible without data (quality data leads to optimized result)

Work-related Factors

- Workload, Job control, Balance.

(contributing to Employee Burnout)

Role clarity, Work-life

Personal Factors

- Age, Gender, personality traits, Emotional intelligence

Organizational Factors

- Organizational culture, leadership style, support systems.

Linear Regression

General Formula:-

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

y - dependent variable (outcome of employee Burnout)

x_1, x_2, \dots, x_n - independent variables (factors which influence dependent variables)

coefficients that give you co-relation b/w Dependent & independent Variable.

β_0 - intercept

ϵ - error term (difference b/w observed & predicted values of y)

Google Co-Lab: online platform where you can implement Machine Learning & Data Science related

(The dataset consist of 9 columns.

- Employee ID
- Date of Joining
- Gender
- Company Type
- WFH Set-up available
- Designation
- Resource allocation
- Mental Fatigue Score
- Burn Rate.

[output column]

} input column

Libraries:
numpy : numerical python (Mathematical Array related computations) to support

why numpy ? → 'n' dimensions (like Matrices, Arrays)

pandas: data Manipulation (dataset from different platforms)

matplotlib, pyplot used for visualization, graphs, animated graphs, charts

seaborn relation bw data with visualization (pictorial rep)

`sklearn.model_selection` : In this module we divide data for testing model (X & Y) after training.

`sklearn.preprocessing` : Converting statistical data into unit variance with the help of standard scalar.

`sklearn.linear_model` - Classification Model : Linear Regression Model. (created using training data-set)

`sklearn.metrics` : to evaluate the performance of our model, metrics i.e. we have; mean_squared_error, mean_absolute_error, r2_score

Pickle } method converter
OS { storing transmitted data from
↓ one format to another
opening a file, reading a file etc.

Data Overview:

1) `data.head()` -

in default case it gives the first '5' entries.

Top '5' entries of the data-set

2) `data.tail(3)` - displays the last 3 entries of a data-set.

3) `data.info()` - it shows the size of each column.

3) `data.describe()` - it gives the statistical information of your data set, (Min, Max, count values... mean, SD)

4) `data.columns.tolist()` - gives you a list of column values in the data set

5) `data.unique()` - No. of unique values under each category / columns it will show you count

6) `data.info()` - Gives the information of your dataset.

7) `data.isnull().sum()` - Count of Null values in each column.

Returns boolean values
whether given column has Null Value or not

8) `data.isnull().sum().values.sum()` - Count of total No. of Null values, in whole data set

9) `data.corr(numeric_only=True)[['Burn Rate'][:-1]]`
- Finds the correlation b/w BurnRate & all other columns having only numerical values.

→ '-1' → it excludes to find correlation among Burn Rate itself

10.) `sns.pairplot(data)`
`plt.show()`

- visualization of correlation of data in terms of the graphical representation.

(correlation between burnout rate & remaining 3 columns) mental fatigue, resource allocation, Observation: to check if designation infer the co-correlation b/w variables

11.) `data = data.dropna()` - Drop off all Null values of our dataframe.
why? because due to the missing values we have some less correlation b/w the models/variables that's why we drop all Null values.

12.) `data.shape` - to verify whether values are dropped off or not.

13.) `data.dtypes` - Analyze data type of each column.

14.) `data = data.drop(['Employee ID', axis=1])`
- dropping the entire employee ID column as it consists of inconsistent data.

To check Employee seniority can contribute much to the data set or Not.

Checking the correlation of Date of Joining with Target Variable.

The answer is, Date of joining to No. of years of work we can find the seniority.

One-hot Encoding, for categorical values

- Based on Number of categories it will be converted into 0 & 1's.
- (Converting categorical features into numerical features)
- dividing the data into input & output.
y - output
x - input
- split the data into training & testing
training size: 70%.
testing size: 30%.

We create a Linear Regression Model where, we feed both the input & output.

Then, we evaluate the model based on error metrics.

- For calculating the errors we will be feeding the machine both Predicted & original outcome so it can generate the error.

Final Observation:

- Based on the evaluation metrics, the Linear Regression model appears to be the best model for Predicting burnout analysis.
- It has the lowest mean squared error, root mean squared error, and mean absolute error, indicating better accuracy and precision in its predictions.
- Additionally, it has highest R-squared score indicating a good fit to the data and explaining a higher & explaining a higher proportion of the variance in the target variable.
- So, we are choosing this model for deployment.

Output:

Linear Regression Model performance Metrics :

Mean Squared Error : 0.0031569779113610717

Root Mean Squared Error: 0.0561869905882231

Mean Absolute Error: 0.04595032644773

R-squared Score : 0.918822674247248.