

# COMPLETE ML, NLP BOOTCAMP MLOPS & DEPLOYMENT

## PYTHON

### Numpy:

- `np.arange(0, 10, 2).reshape(5, 1)` →  $\begin{bmatrix} [0] \\ [2] \\ [4] \\ [6] \\ [8] \end{bmatrix}$
- `np.ones((3, 4))` →  $\begin{bmatrix} [1, 1, 1, 1] \\ [1, 1, 1, 1] \\ [1, 1, 1, 1] \end{bmatrix}$
- `np.eye(3)` → forms identity matrix  $\begin{bmatrix} [1, 0, 0] \\ [0, 1, 0] \\ [0, 0, 1] \end{bmatrix}$
- `arr = np.array([1, 2, 3, 4], [5, 6, 7, 8], [9, 10, 11, 12])`  
`print(arr[0:2, 2:3])` →  $\begin{bmatrix} [3, 4] \\ [7, 8] \end{bmatrix}$
- Numpy array of shape (6, 6) with values = 1 to 36. → `np.arange(1, 37).reshape(6, 6)`  
`array = print(arr)`
- Numpy array of shape (5, 5) filled with random integers. Use fancy indexing to extract the elements at the corners of the array.  
`array = np.random.randint(1, 21, size=(5, 5))`  
`print(array)`  
`corners = array[[0, 0, -1, -1], [0, -1, 0, -1]]`  
`print(corners)`

} circled ele. is = output.

### Pandas:

- Create a series from dictionary: `data = {'a': 1, 'b': 2}` → `pd.Series(data)`  
 $\underbrace{\quad}_{\text{1D array}}$        $\left. \begin{array}{l} \text{a} : 1 \\ \text{b} : 2 \end{array} \right\}$

For 2D array → Dataframe

### Data aggregating & grouping:

`grouped_mean = df.groupby('Product')['Value'].sum()`

### Merging & joining Dataframe

```
df1 = pd.DataFrame({'key': ['A', 'B', 'C'], 'value1': [1, 2, 3]})  

df2 = pd.DataFrame({'key': ['A', 'B', 'D'], 'value2': [1, 5, 6]})  

pd.merge(df1, df2, on="key", how="inner")  

pd.merge(df1, df2, on="key", how="outer")
```

### Reading Data from Different Sources

→ `from io import StringIO`

`!pip install html`

`!pip install html5lib`

• `df = pd.read_html(url)`

`!pip install beautifulsoup4`

• `df_excel = pd.read_excel('data.xlsx')`

## Seaborn:

```

import seaborn as sns
tips = sns.load_dataset('tips')
sns.scatterplot(x='total_bill', y='tip', data=tips)
sns.lineplot(x='size', y='total_bill', data=tips)
sns.barplot(x='day', y='total_bill', data=tips)
sns.boxplot()
sns.violinplot()
sns.histplot(tips['total_bill'], bins=10, kde=True)
sns.kdeplot(tips['total_bill'], fit=True)
sns.pairplot(tips)
corr = tips[['total_bill', 'tip', 'size']].corr()
sns.heatmap(corr, annot=True, cmap='coolwarm')
    
```

↳ Shows relation b/w all Variables

## SQLite:

It's a lightweight, self-contained, serverless & zero-configuration SQL database engine.

↳ import sqlite3

- Connect to SQLite DB: `connection = sqlite3.connect('example.db')`

- After this, `example.db` will be created in current directory

Now to execute queries on the DB;

• `cursor = connection.cursor()`

- Create a table

```

cursor.execute('''
Create Table If Not Exists employees (
    id Integer primary key,
    name Text not null,
    age Integer,
    department Text
)''')
connection.commit()
    
```

- Inserting the data

```

cursor.execute('''
insert into employees(
    name, age, department)
values ('Ankit', 32, 'SWE')
''')
connection.commit()
    
```

- Querying the data

```

cursor = execute('select * from employees')
rows = cursor.fetchall()
    
```

```

for row in rows:
    print(row)
    
```

- Insert multiple rows

```

sales_data = [('2023-01-01', 'Product 1', 100), ('2023-01-02', 'Product 2', 200)]
cursor.executemany('''
insert into sales(date, product, sales, region)
values (?, ?, ?, ?)
''', sales_data)
connection.commit()
    
```

## File Operations - Read and Write Files

- Read a whole file:

```

with open('example.txt', 'r') as file:
    content = file.read()
    print(content.strip())
    
```

- Write a file with overwriting:

```

with open('example.txt', 'w') as file:
    file.write('I am kind \n')
    
```

- Write a file without overwriting:

```

with open('example.txt', 'a') as file:
    file.write('Appending this line \n')
    
```

- Writing a list of lines to a file: ↳ `file.writelines(lines)`

I  
(list)

→ The 'wt' mode in Python is used to open a file for both reading and writing. If the file does not exist, it will be created. If the file exists, the content in file is overwritten.

with open('example.txt', 'wt') as file:

file.write("Hello World \n")

file.seek(0) → moving cursor to the beginning

content = file.read()

print(content)

## Working with File Paths:

- import os

`cwd = os.getcwd()` → gets current working directory

- `items = os.listdir('.') → Listing files & directories`

`print(items)`

(let path = 'example.txt')

- `os.path.isfile(path) → checks if 'path' is a file`

`os.path.isdir(path) → checks if 'path' is a directory`

## Exception Handling: try, catch & finally

def divide(a,b):

try:

`result = a/b`

except ZeroDivisionError as e:

`print(f"Error {e}")`

`result = None`

finally:

`print('Execution Complete!')`

`return result`

`print(divide(10,0))`

↳ This block always runs independent of try or catch.

# #OOPS

## Inheritance:

class Car:

```
def __init__(self, doors, engine_type):
    self.doors = doors
    self.engine_type = engine_type

def drive(self):
    print(f"Engine type: {self.engine_type}. car.")
```

car1 = Car(5, "petrol")

car1.drive()

class Tesla(Car):

```
def __init__(self, doors, engine_type, self_driving):
    super().__init__(doors, engine_type)
    self.is_self_driving = self_driving

def self_driving(self):
    print(f"Supports self driving: {self.is_self_driving}")
```

t1 = Tesla(5, "electric", True)

t1.self\_driving()

t1.drive()

↳

petrol car	electric car
Supports: True	Supports: False

## In Multiple Inheritance:

class Animal:

// constructor ~ (self, name)

def speak(self)

class Pet:

// constructor ~ (self, owner)

class Dog(Animal, Pet):

```
def __init__(self, name, owner):
    Animal.__init__(self, name)
    Pet.__init__(self, owner)
```

## Polymerism

### Method Overriding:

class Veh:

```
def start(self):
    print("V")
```

class Car(Veh):

```
def start(self):
    print("C")
    super().start()
```

Car = Car()

car.start()

↳

C	V
---	---

ex: polymor. with func.'s & methods.

class Shape:

```
def area(self):
    return "Area is"
```

class Rectangle(Shape):

```
def __init__(self, width, height):
    self.width = width
    self.height = height
```

def area(self):

```
return self.width *
       self.height
```

def print\_area(shape):

```
print(f"Area = {shape.area()}"")
```

rectangle = Rectangle(4, 3)

circle = Circle(3)

print\_area(rectangle)

print\_area(circle)

class Circle(Shape):

```
def __init__(self, rad):
    self.rad = rad
```

def area(self):

```
return 3.14 * rad * rad
```

ex: Operator Overloading: allows to define the behaviour of operators for custom objects.

- common operator overloading magic methods:

- \_\_add\_\_(self, other) : for +
- \_\_sub\_\_(self, other) : for -
- \_\_mul\_\_ : for \*
- \_\_truediv\_\_ : for division
- \_\_eq\_\_ : for checking ==

ex: class Vector:

```
def __init__(self, x, y):
    self.x = x
    self.y = y
```

```
v1 = Vector(2, 3)
v2 = Vector(4, 5)
print(v1 + v2)
```

def \_\_add\_\_(self, other):

```
return Vector(self.x + other.x,
              self.y + other.y)
```

• Abstract Class: It is a class that is used as a base class for defining abstract classes, which cannot be instantiated directly.

(It's blueprint for other classes. It can't be directly used to create objects.)

from abc import ABC, abstractmethod

class Vehicle(ABC):

\_\_@abstractmethod

def start(self):

pass

class Car(Vehicle):

def start(self):

return "Car"

class Bike(Vehicle):

def start(self):

return "Bike"

def Veh(Vehicle):

print(Vehicle.start())

Car = Car()

Veh(Car)

↳

Car	Veh(Car)
-----	----------

The 'Vehicle' class can't be used on its own to create objects. It only defines that every 'Vehicle' must have 'start' method but doesn't specify what it does.

## Encapsulation:

```
class Person:
    def __init__(self, name, age, gender):
        self.__name = name → private access Modifier
        self.__age = age → protected access
        self.gender = gender → public access
```

```
def get_name(person):
    return person.gender
person = Person("Ankit", 19, "M")
get_name(person) → M
```

dir(person):  
lists all attributes & methods available for the 'person' object

Abstraction: → hides complex implementation details

Magic Methods: predefined methods in Python that one can override to change behaviour of the objects.

Ex: → `__init__`: initializes a new instance of a class  
`__str__`: returns a string representation of an object  
`__repr__`: returns an official string representation of an object  
`__len__`: returns the length of an object

```
class Person:
    def __init__(self, name, age):
        self.name = name
        self.age = age

    def __str__(self):
        return f"{self.name}, {self.age} years"

    def __repr__(self):
        return f"Person(name={self.name})"
```

```
person = Person("Ankit", 19)
print(person)
print(repr(person))
↓
obj:
Ankit, 19 years
Person(name=Ankit)
```

## Custom Exception:

```
class Error(Exception):
    pass

class dobException(Error):
    pass

year = int(input("Enter your dob:"))
age = year - 2024

try:
    if age <= 30 and age >= 20:
        print("Yes")
    else:
        raise dobException()
except dobException:
    print("No")
```

(In java/C++ → try-catch)  
 In python → try-except

(when we are raising an exception, we also need to catch it, so writing 'except' block.)

# Iterators: Allows efficient looping & memory management.

```
ex: m = [1, 2, 3, 4]
iterator = iter(m)
next(iterator) → outputs next ele. each time we run this code block
                (Gives "StopIteration" error at end)
                , uses lazy loading technique.
```

# Generators: are a simpler way to create iterators

```
def square(n):
    for i in range(3):
        yield i**2 → creates a local variable & it stores the value
a = square(3)
next(a)
```

# Decorators: Allows us to modify the behaviour of a func. or class method.

↳ 3 features: function copy, closures, decorators.

- `function_copy`: (we can copy 1 func. to other)
 

```
def welcome():
    return "Hi"
wel = welcome
print(wel())
# output: Hi
```
- `Closures function`: a function inside a function
 

```
def main_wel(func, lst):
    def sub_wel():
        print("Hi")
        print(func(lst))
    return sub_wel
main_wel(len, [1, 2, 3])
# output: Hi
#           3
```
- `Decorators`: allows the wrapper to accept any no. of positional arguments
 

```
*args: allows the wrapper to accept any no. of positional arguments
```

`def my_decorator(func):`  
 `def wrapper():
 print("Before")
 func()
 print("After")
 return wrapper`

`@my_decorator`  
`def say_hello():
 print("Hi")`

`say_hello()` → obj: Before  
 Hi  
 After

`**kwargs`: allows the wrapper to accept any no. of keyword arguments

`def example(*args, **kwargs):`

`print("Args:", args)`

`print("Kwargs:", kwargs)`

`example(1, 2, 3, name="AV", age=19)`

`Output: Args: (1, 2, 3)`

`Kwargs: {'name': 'AV', 'age': 19}`

Python Logging: built-in logging module offers a flexible framework for emitting log messages from Python programs.

Several log levels:

- DEBUG: detailed info., typically of interest only when diagnosing problems.
- INFO: confirmation that things are working as expected.
- WARNING: something unexpected happened.
- ERROR: A more serious issue, the software has not been able to perform some func.
- CRITICAL: Very serious error.

import logging

```
logging.basicConfig(level=logging.DEBUG)
logging.debug("debug message")
```

→ DEBUG: debug message

→ Program: a sequence of instructions written in programming language

→ Process: an instance of a program that is being executed.

→ Thread: unit of execution within a process

Multiprocessing: → used when we do "parallel execution"

```
if __name__ == "__main__":
    import multiprocessing
    p1 = multiprocessing.Process(target=square_numbers)
    p2 = "" # another process
    t = time.time()
    p1.start() → Start the process
    p2.start() → wait for process to complete
    p2.join()
    print(time.time() - t) → prints time to complete
```

• We can also use "Process Pool Executor"

Multi threading:

```
if __name__ == "__main__":
    import threading
    import time
    t1 = threading.Thread(target=print_numbers)
    t2 = "" # another thread
    t = time.time()
    t1.start()
    t2.start()
    t1.join()
    t2.join()
    print(time.time() - t)
```

• We can also use "Thread Pool Executor"

Memory Management:

• Reference Counting: is the primary method Python uses to manage memory.

Each object in python maintains a count of references pointing to it. When the reference count = 0, the memory occupied by the object is deallocated.

```
import sys
a = []
print(sys.getrefcount(a))
```

one reference from a

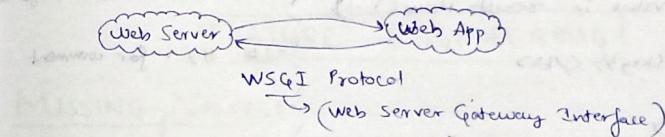
one from getrefcount()

• Garbage Collection: Python includes a cyclic garbage collector to handle reference cycles. Reference cycles occur when objects reference each other, preventing their reference counts from reaching zero.

• Best practices:

- 1) Use local variables
- 2) Avoid circular reference
- 3) Use generators
- 4) Explicitly delete objects
- 5) Profile Memory Usage

**FLASK** → web framework



from flask import Flask, render\_template, request

```
app = Flask(__name__)
@app.route("/") → will be my WSGI application.
def welcome():
    return "<html><h1> Welcome </h1></html>"
```

responsible in redirecting to another HTML page. It looks for a folder 'templates' & it will redirect to the html pages present in it.

```
@app.route("/index", methods=['GET'])
def index():
    return render_template('about.html')
```

```
@app.route("/form", methods=['GET', 'POST'])
def form():
    if request.method == 'POST':
```

```
        name = request.form['name']
        return f"Hello {name}!"
    return render_template('form.html')
```

```
if __name__ == "__main__":
    app.run(debug=True)
```

→ Results in running entire flask app.

• To run above app in terminal

python <filename>

'debug=True' → results in dynamic web pages if any changes in code is done. It restarts the server automatically.

We can also define : (POST, PUT, DELETE) operations in Flask

dynamic url

```
@app.route('/successres/')
def successres(score):
    res = ''
    if score >= 50:
        res = "PASSED"
    else:
        res = "FAILED"

    exp = {'score': score, 'res': res}
    return render_template('result.html', results=exp)
```

Now in result.html:  $\Rightarrow$  Jinja2 template: default template engine used by flask

(allows python like expressions to embed into HTML for creating dynamic web content.)

```
<html>
<h2> Final Results </h2>
<body>
<table border=2>
    {%- for key,value in results.items() -%}
    <tr>
        <td> {{key}} </td>
    </tr>
    {%- end for -%}
</table>
</body>
</html>
```

• {{ }} - expression to print output in html  
 • {%- %} - condition for loop  
 • # - for comment.

## STREAMLIT

It's an open-source app framework for ML & DS projects. It allows us to create beautiful web applications with simple python scripts.

To run it  $\rightarrow$  open terminal  $\rightarrow$  "streamlit run <filename>"

```
import streamlit as st
import pandas as pd
import numpy as np

st.title("Hello")
df = pd.DataFrame(
    np.random.randn(20, 3), columns=['a', 'b', 'c']
)
st.line_chart(df)
```

① Merge 2 lists into dictionary.

```
A def func(key, values):
    if len(key) != len(values):
        return False
    d = dict(zip(key, values))
    return d
```

↳ makes into key-val pairs

② Merge multiple dictionaries into 1

```
A def func(d1, d2, d3):
    d = d1 | d2 | d3
    return d
```

$\Rightarrow$  If same key present in other dictionaries, the value of rightmost dictionary in above will overwrite the prev. value.

③ Count frequency of words in sentence.

```
A def func(sentence):
    sentence = sentence.lower()
    word = sentence.split()
    word_count = {}
    for word in words:
        if word in word_count:
            word_count[word] += 1
        else:
            word_count[word] = 1
    return word_count
```

④ Check if tuple is palindromic

```
A def func(tup):
    return tup == tup[::-1]
```

⑤ Merge dictionaries with common keys  $\rightarrow$  add their value, if key  $\rightarrow$  common

```
A def func(dict):
    ans = {}
    for d in dict:
        for key, value in d.items():
            if key in ans:
                ans[key] += value
            else:
                ans[key] = value
    return ans
```

## FEATURE ENGINEERING

### MISSING VALUES:

1) Missing Completely at Random (MCAR): Missing values are completely random & not related to any observed or unobserved data. The probability of a value being missing is the same for all observations.

2) Missing at Random (MAR): Missing values are related to observed data but not to the missing value themselves. Given the observed data, the probability of a value being missing can be estimated.

3) Missing Not at Random (MNAR): The missing data is directly related to the value that is missing.

### Handling Missing Values:

1) Mean Value Imputation: filling with 'mean'

```
ex: df['Age-mean'] = df['age'].fillna(df['age'].mean())
```

2) Median Value Imputation: {If we have outliers in the dataset}

```
ex: df['Age-median'] = df['age'].fillna(df['age'].median())
```

3) Mode Imputation: {Categorical Values}

↳ filling with 'mode' of that column.

### Handling Imbalanced Dataset:

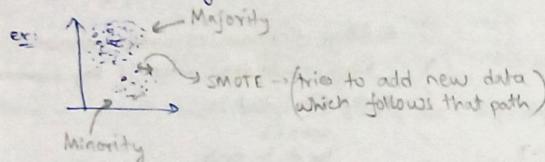
1) UpSampling: Increasing the no. of instances in the minority class to balance the dataset.

// refer example in "notebook" → used → from sklearn.utils import resample

- Basically, it duplicates the same data to match the 'no.' of majority class.
- 2) DownSampling: Reducing the no. of instances in the majority class to balance the dataset.  
(Generally not preferred as there's loss of data.)

### 3) SMOTE: {Synthetic Minority Oversampling Technique}

It generates synthetic instances of minority classes by interpolating between existing instances.



from imblearn.over\_sampling  
import SMOTE  
refer example in "notebook."

\* lst = [45, 32, ..., 74] # a list

import seaborn as sns  
sns.boxplot(lst)

My boxplot will automatically generate  
"5 number summary" and draw a box  
(HANDLING OUTLIERS)

## DATA ENCODING:

Means that we will be converting the categorical features into some numerical values, so that the model will be able to understand for a specific purpose. We have 3 types of encoding:

### 1) Nominal / One-Hot Encoding:

Here each category is represented as binary vector where each bit corresponds to a unique category. Ex: If we have a cat. var. "color" with 3 possible values (R, G, B), we can represent it using one hot encoding as follows:

↳ Red : [1, 0, 0]  
↳ Green : [0, 1, 0]  
↳ Blue : [0, 0, 1]

Disadvantage: If we have categorical variable with many features (like 100), then it will create 100 new columns for this, which is not recommended.

ex: from sklearn.preprocessing import OneHotEncoder

df = pd.DataFrame({ 'color': ['red', 'blue', 'green', 'green', 'red']})

encoder = OneHotEncoder() # creating an instance of OHE

encoded = encoder.fit\_transform(df[['color']]).toarray() # 2D → converts into a sparse matrix

encoder\_df = pd.DataFrame(encoded, columns=encoder.get\_feature\_names\_out())

encoder\_df →

	color_blue	color_green	color_red
0	0	0	1
1	1	0	0
2	0	1	0
3	0	1	0

# we can also concat.

pd.concat([df, encoder\_df], axis=1)

### 2) Label Encoding:

It involves assigning a unique numerical label to each category in the variable. The labels are usually assigned in alphabetical order or based on the frequency of the categories.

ex: Color → Red: 1 } or from sklearn.preprocessing import LabelEncoder  
Green: 2 }      lbl = LabelEncoder()  
Blue: 3 }      lbl.fit\_transform(df[['color']])

Disadvantage: Here red=2, blue=0, green=1. The model might start thinking that red have higher value/priority than blue. This should not happen because these are nominal values.

But yes, if there's a situation where we need to assign ranks to these particular values to this category, we can definitely do it with the help of ordinal encoding.

### 3) Ordinal Encoding:

In this, each category is assigned a numerical value based on its position in the order.

ex: education level → High School : 1 } ← represented via ordinal encoding  
College : 2  
Graduate : 3  
Post-Graduate : 4

ex: from sklearn.preprocessing import OrdinalEncoder  
df = pd.DataFrame({ 'size': ['small', 'medium', 'large', 'medium', 'small', 'large']})  
encoder = OrdinalEncoder(categories = [['small', 'medium', 'large']])  
encoder.fit\_transform(df[['size']])

# New data  
encoder.transform(['small']) → array([0])

### 4) Target Guided Ordinal Encoding:

It's used to encode categorical variables based on their relationship with the target variable.

This technique is useful when we have a categorical variable with a large no. of unique categories, & we want to use this variable as a feature in our machine learning model.

Here, we replace each category in the categorical variable with a numerical value based on the mean or median of the target variable for that category.

ex: df = pd.DataFrame({ 'city': ['A', 'B', 'C', 'D', 'A', 'C'], 'price': [10, 20, 30, 25, 23, 40]})  
mean\_price = df.groupby('city')[['price']].mean().to\_dict()  
mean\_price → dict: {'A': 19.5, 'B': 20.0, 'C': 23.0, 'D': 20.0}  
df['city-encoded'] = df['city'].map(mean\_price)

	Price	city_encoded
0	10	17.5
1	20	20.0
2	30	35.0

## Practicing EDA:

- 1) Red Wine Data:
  - .head() → summary of dataset
  - .info() → descriptive summary of dataset
  - checking unique & missing values. finding duplicate records & dropping them
  - Using 'sns' visualizations to bring out conclusions → like: pairplot, histplot, scatterplot, etc.

- 2) Flight Price Prediction:
  - .head(), .info(), .describe()

→ splitting columns → Date of Journey → Day | Month | Year  
 24/03/2019      24    03    2019

→ Also checking the dtype of all columns. Ensure numeric columns have dtype as 'int' or 'float' & not 'obj'

→ Replacing 'nan', empty blocks with 'mode' of that column  
 → Using OHE to change categorical values into numerical values

- 3) Google Play Store Dataset:

→ .shape, .info(), .describe(), .isnull().sum()

# In EDA, we just need observation.

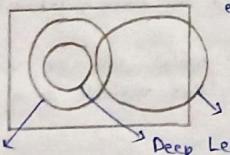
→ converting col. values into meaningful data → ex: In reviews → 3M+  
 (changing to int)

→ Removing duplicate records.

→ Plotting different parameters & noting observations  
 (like: top 10 apps used)

## INTRODUCTION TO MACHINE LEARNING

AI = to create an application which can perform its own task without any human intervention  
 ex: Netflix recommendation system



Data Science

ML

Deep Learning (Multi-Layered NN) = mimic the human brain

(It provides stats tools to analyze, visualize, predict & forecast the data sometimes)

Machine Learning Technique

Supervised ML

Unsupervised ML

Reinforcement Learning

- Supervised ML:
  - Regression: predicts continuous values. {continuous}
  - Classification: predicts discrete class labels. {categorical}
  - Dependent on output feature
  - binary classification = only 2 categories
  - Multiclass classification } > 2 categories

- Linear Regression
- Ridge & Lasso
- ElasticNet
- Logistic Regression (Classification)
- Decision Tree
- Random Forest
- Ada Boost
- XGBoost
- both classification & regression.

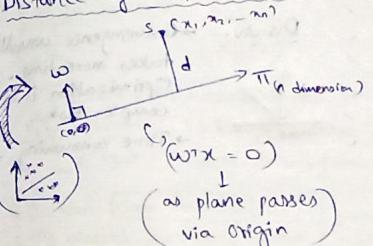
- Unsupervised ML (we don't have o/p feature label)

We create clusters to group similar kind of data

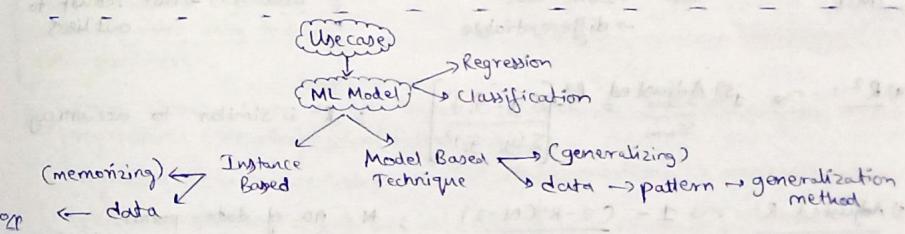
- k Means
- Hierarchical Mean
- DBScan Clustering

Refer "Eqn. of a Line, Plane, Hyperplane" → PREV. CHAPTER

Distance of a point from plane:



$$d = \frac{w^T.s}{\|w\|} ; w^T.s = \|w\| \|s\| \cos \theta$$



### Usual ML

- Prepare the data for model training
- Train model from training data to estimate model parameters i.e. discover patterns
- Store the model in suitable form
- Generalize the rules in form of model, even before scoring instance is seen
- Predict for unseen scoring instance using model
- Can throw away training data after Model training
- Requires a known model form
- Scoring models generally require less storage
- Scoring for new instance is generally faster.

### Instance Based Learning

- Prepare the data for model training. No diff here
- Do not train model - pattern discovery postponed until scoring query received
- There is no model to store.
- No generalization before scoring. Only generalize for each scoring instance individually as and when seen.
- Predict for unseen scoring instance using training data directly
- Training data must be kept since each query uses part/full set of training observations
- May not have explicit model form.
- Storing training data generally requires more storage.
- Scoring for new instance may be slow.

# UNDERSTANDING LINEAR REGRESSION INDEPTH INTUITION AND PRACTICALS

Refer "Application of Linear Algebra, Stats & Differential Calculus" → PREV. CHAPTER  
Convergence Algorithm:

$$\text{Repeat until convergence} \Rightarrow \theta_j = \theta_j - \alpha \left( \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \right); \quad J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\theta_0 = \theta_0 - \alpha \left( \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \right) \quad \theta_1 = \theta_1 - \alpha \left( \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_i \right)$$

// Refer Prev. Chapter

Cost Function → Performance Metrics: → I/W Que.: Check usage of all performance metrics

1) MSE:  
 $\downarrow$   
 $\sum_{i=1}^n (y_i - \hat{y}_i)^2 / n$   
 Adv.: → Differentiable  
 → It has 1 local & 1 global minima.  
 → Converges faster  
 Disadv.: → Not robust to outliers  
 → It is not in same unit.

2) MAE: (Mean Absolute Error)  
 $\downarrow$   
 $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$   
 Adv.: → Robust to outliers  
 → It will be in the same unit.  
 Disadv.: → Convergence usually takes more time.  
 Optimization is a complex task.  
 → Time consuming.

3) RMSE: (Root Mean Squared Error)  
 $\downarrow$   
 $\sqrt{MSE}$   
 Adv.: → same unit  
 → differentiable  
 Disadv.: → Not robust to outliers

4) R<sup>2</sup>:  $\sim 1 - \frac{SS_{res}}{SS_{tot}}$   $= 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$  ; R<sup>2</sup> is similar to accuracy.

5) Adjusted R<sup>2</sup>:  $\sim 1 - \frac{(1-R^2)(N-1)}{N-p-1}$ ; N = no. of data points  
 p = no. of independent features  
 // Refer Prev. Chapter

Linear Regression Using OLS (Ordinary Least Square):

$h_\theta(x) = \beta_0 + \beta_1 x_i$  → OLS: calculates  $\beta_0$  &  $\beta_1$   
 (main aim: reduce the error)

$\beta_0 = \bar{y} - \beta_1 \bar{x}$  } Intercept

$\beta_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})}$  } Coefficient

(for derivation → refer pdf)  
 Basically, cost func. is used to differentiated w.r.t  $\beta_0$  &  $\beta_1$

Exploratory Data Analysis (EDA) using Python:

Ex: Refer Notebook for output

df = pd.read\_csv('height-weight.csv')  
 df.head() → exploring dataset  
 plt.scatter(df['Weight'], df['Height']) → scatterplot for visualization  
 df.corr() → Checking correlation b/w variables  
 import seaborn as sns  
 sns.pairplot(df)

I M P L E M E N T A T I O N

# Independent & Dependent features  
 $x = df[['Weight']]$  → To create independent features, ensure it's in the form of a Dataframe or 2D array.  
 $y = df['Height']$  → Dependent / Target variable can be series or 1D array.

L # Train-test-split  
 from sklearn.model\_selection import train\_test\_split  
 x\_train, x\_test, y\_train, y\_test = train\_test\_split(x, y, test\_size=0.25, random\_state=42)

R # Standardization  
 from sklearn.preprocessing import StandardScaler  
 We apply this because the i/p values might be higher & we need to draw a gradient descent, so higher the value, more the time it will take. So we apply Z-score & convert it into  $\mu=0, \sigma=1$ , so that it takes less time.

G (Z-score)  $\frac{(x_i - \mu)}{\sigma}$   
 Z-score transformation occurs in all values of x-train

S Scaler = StandardScaler()  
 x\_train = Scaler.fit\_transform(x\_train)  
 x\_test = Scaler.transform(x\_test) } 'transform' represents we will use same  $\mu$  &  $\sigma$  as the one used in training dataset.

N To avoid [data leakage], as my model shouldn't know anything about test data, so we are using only training data parameters. WHY?

① from sklearn.linear\_model import LinearRegression → Applying ML Model  
 regression = LinearRegression(n\_jobs=-1)  
 regression.fit(x\_train, y\_train) } 'n\_jobs=-1' = use all processor available in the system.  
 print(regression.coef\_) →  $\beta_1$  → means 1 unit movement in x equals ' $\beta_1$ ' movement in y  
 print(regression.intercept\_) →  $\beta_0$   
 plt.scatter(x\_train, y\_train)  
 plt.plot(x\_train, regression.predict(x\_train)) → Visualizing  
 # Prediction of test data  
 y\_pred = regression.predict(x\_test) } height or Predicted = intercept + coef\*(weight)  
 from sklearn.metrics import mean\_absolute\_error, mean\_squared\_error  
 mse = mean\_squared\_error(y\_test, y\_pred)  
 mae = mean\_absolute\_error(y\_test, y\_pred)  
 rmse = np.sqrt(mse)  
 print(mse, mae, rmse) } Performance metrics

② from sklearn.metrics import r2\_score  
 score = r2\_score(y\_test, y\_pred)  
 print(score)  
 $1 - (1 - \text{score}) * (\text{len}(y\_test) - 1) / (\text{len}(y\_test) - x\_test.\text{shape}[1] - 1)$  } Adjusted R<sup>2</sup> score

- 0: from sklearn.model\_selection import train\_test\_split  
 $x_{-train}, x_{-test}, y_{-train}, y_{-test} = \text{train\_test\_split}(x, y, \text{test\_size}=0.25, \text{random\_state}=42)$
- 1: from sklearn.preprocessing import StandardScaler  
scaler = StandardScaler()  
 $x_{-train} = \text{scaler}.\text{fit\_transform}(x_{-train})$   
 $x_{-test} = \text{scaler}.\text{transform}(x_{-test})$
- 1.1
- from sklearn.linear\_model import LinearRegression  
regression = LinearRegression(n\_jobs=-1)  
regression.fit(x\_train, y\_train)
- 2: from sklearn.metrics import mean\_absolute\_error, mean\_squared\_error  
mse = mean\_squared\_error(y\_test, y\_pred)  
mae = mean\_absolute\_error(y\_test, y\_pred)  
rmse = np.sqrt(mse)  
print(mse, mae, rmse)
- from sklearn.metrics import r2\_score  
score = r2\_score(y\_test, y\_pred)
- (adjusted R2-score)
- print(score)  
 $1 - (1 - \text{score}) * (\frac{\text{len}(y_{-test}) - 1}{\text{len}(y_{-test}) - x_{-test}.\text{shape}[1] - 1})$
- 3: import statsmodels.api as sm  
model = sm.OLS(y\_train, x\_train).fit()  
prediction = model.predict(x\_test)  
print(prediction)  
model.summary  
regression.predict(scaler.transform([[72]]))
- 4: from sklearn.preprocessing import PolynomialFeatures  
poly = PolynomialFeatures(degree=2, include\_bias=True)  
x\_train\_poly = poly.fit\_transform(x\_train)  
x\_test\_poly = poly.transform(x\_test)  
regression = LinearRegression()  
regression.fit(x\_train\_poly, y\_train)  
y\_pred = regression.predict(x\_test\_poly)  
score = r2\_score(y\_test, y\_pred)  
print(score)
- 5: linreg = LinearRegression()  
linreg.fit(x\_train\_scaled, y\_train)  
y\_pred = linreg.predict(x\_test\_scaled)  
mae = mean\_absolute\_error(y\_test, y\_pred)  
score = r2\_score(y\_test, y\_pred)  
plt.scatter(y\_test, y\_pred)
- 6: from sklearn.metrics import accuracy\_score, classification\_report, confusion\_matrix  
print(accuracy\_score(y\_test, y\_pred))  
print(classification\_report(y\_test, y\_pred))  
print(confusion\_matrix(y\_test, y\_pred))
- 7: from sklearn.linear\_model import LogisticRegression  
logistic = LogisticRegression()  
logistic.fit(x\_train, y\_train)  
y\_pred = logistic.predict(x\_test)  
print(y\_pred)  
logistic.predict\_proba(x\_test)

## SHORTCUTS

8: Model = LogisticRegression()

params = {  
 'penalty': ['l1', 'l2', ''],  
 'C': [100, ...],  
 'solver': ['sag', ...]}

9: grid = GridSearchCV(estimator=Model, param\_grid=params, scoring='accuracy', cv=cv,  
grid.fit(x-train, y-train)  
grid.best\_params\_, grid.best\_score\_

y-pred = grid.predict(x-test)

10: param\_grid = {'C': [0.1, ...],  
 'gamma': [1, ...],  
 'kernel': ['rbf', ...]}

grid = GridSearchCV(SVC), param\_grid=param\_grid(), refit=True, cv=5, verbose=3)

grid.best\_params\_

y-pred2 = grid.predict(x-test)

11: Refer  $\Rightarrow$  "Random Forest Machine Learning"

12: ↗ ↘ 13: ; 14: Refer  $\Rightarrow$  "NLP For Machine Learning"  $\rightarrow$  (BOW + N-grams)

## # OLS Linear Regression

```
import statsmodels.api as sm
model = sm.OLS(y-train, x-train).fit()
prediction = model.predict(x-test)
print(prediction) → prints predicted values of x-test.
```

Model.summary() → Gives full summary of model (like coef, std. err, etc)

# Prediction for new data

regression.predict(scalar.transform([[72]])) → we need to standardize (then predict)

# Here all 'plt' diagrams styling written in brief → refer notebook for detailing.  
Ex. (Multiple Linear Regression) ; Refer Notebook for output

```
df = pd.read_csv('economic_index.csv')
```

df.head() → drop unnecessary columns

df.drop(columns = ['Unnamed: 0', 'year', 'month'], axis=1, inplace=True)

df.isnull().sum() → check null values

# Visualization

import seaborn as sns

sns.pairplot(df)

df.corr()

# Get independent & dependent features

x = df.iloc[:, :-1] # x = df[['interest\_rate', 'unemployment\_rate']]

y = df.iloc[:, -1] → means take all rows ':', except last row ':-1'

x.head() → means take the last column = 'index\_price'

y

from sklearn.model\_selection import train\_test\_split

x\_train, x-test, y-train, y-test = train\_test\_split(x, y, test\_size=0.25, random\_state=42)

import seaborn as sns

sns.regplot(x=df['interest\_rate'], y=df['index\_price'])

(→ Plots data & regression linear model fit. →

→ Works only for 2 features

sns.regplot(x=df['index\_price'], y=df['unemployment\_rate'])

①

from sklearn.model\_selection import cross\_val\_score

validation\_score = cross\_val\_score(regression, x-train, y-train,

scoring='neg\_mean\_squared\_error', cv=3)

np.mean(validation\_score) → (avg. of '3' tests) = (CV=3)

# Prediction

y-pred = regression.predict(x-test)

y-pred

②

residuals = y-test - y-pred

sns.displot(residuals, kind='kde')

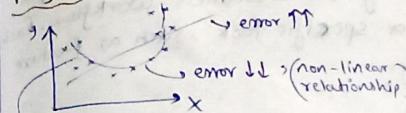
plt.scatter(y-pred, residuals)

↓

Point are uniformly distributed, so we can conclude that there is no issue in our process

③

## Polynomial Regression:



Simple LR : polynomial degree = 0 :  $h_0(x) = \beta_0 \cdot x^0 \Rightarrow$  constant value  
(only 1 feature) → " " " = 1 :  $h_1(x) = \beta_0 \cdot x^0 + \beta_1 \cdot x^1 \Rightarrow$  Simple LR  
→ " " " = 2 :  $h_2(x) = \beta_0 \cdot x^0 + \beta_1 \cdot x^1 + \beta_2 \cdot x^2$   
→ " " " = n :  $h_n(x) = \beta_0 \cdot x^0 + \beta_1 \cdot x^1 + \beta_2 \cdot x^2 + \dots + \beta_n \cdot x^n$

(If 2 independent features)-  
degree = 1 →  $h_1(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$   
degree = 2 →  $h_2(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2$

As degree increases, the curve gets more twisted & the best fit line will come & eventually overfit. We should select a degree value which should not overfit/underfit but generalize the model

Ex. (Polynomial Regression) ; Refer Notebook for output

# creating dataset

x = 6 \* np.random.rand(100, 1) - 3

y = 0.5 \* x \*\* 2 + 1.5 \* x + 2 + np.random.randn(100, 1)

plt.scatter(x, y)

①

# We can scale the values, but they're already small, so it's not necessary

from sklearn.metrics import r2\_score  
score = r2\_score(y-test, regression.predict(x-test))  
print(score)

plt.plot(x-train, regression.predict(x-train))

# Applying Polynomial transformation

from sklearn.preprocessing import PolynomialFeatures

poly = PolynomialFeatures(degree=2, include\_bias=True)

x-train-poly = poly.fit\_transform(x-train)

x-test-poly = poly.transform(x-test)

x-train-poly, x-test-poly

regression = LinearRegression()

regression.fit(x-train-poly, y-train)

y-pred = regression.predict(x-test-poly)

score = r2\_score(y-test, y-pred)

print(score)

print(regression.coef\_, regression.intercept\_)

plt.scatter(x-train, regression.predict(x-train-poly))

plt.scatter(x-train, y-train)

④ → with "degree = 3"

# Prediction of new data set

x-new = np.linspace(-3, 3, 200).reshape(200, 1)

x-new-poly = poly.transform(x-new)

y-new = regression.predict(x-new-poly)

plt.plot(x-new, y-new, label="New Predictions")

plt.plot(x-train, y-train, label="Training Points")

plt.plot(x-test, y-test, label="Testing Points")

(2 up, 2 up features)

$h_0(x) = \beta_0 + \beta_1 x \rightarrow$  Simple linear Regression

$h_3(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 \rightarrow$  Multiple LR

Refer diagram in notebook

Pipelining: It's a process in DSML where multiple sequential steps are assembled into a single, streamlined workflow. Each step in the pipeline typically performs a specific task - such as feature extraction, preprocessing, or modeling.

Ex: from sklearn.pipeline import Pipeline

def poly\_regression(degree):

x\_new = np.linspace(-3, 3, 200).reshape(200, 1)

poly\_features = PolynomialFeatures(degree=degree, include\_bias=True)

lin\_reg = LinearRegression()

poly\_regression = Pipeline([

("poly-features", poly\_features),

("lin-reg", lin\_reg)

)

poly\_regression.fit(x\_train, y\_train) # Polynomial & fit of LR

y\_pred\_new = poly\_regression.predict(x\_new)

# Plotting prediction line

plt.plot(x\_new, y\_pred\_new, 'r', label="Degree " + str(degree), linewidth=3)

plt.plot(x\_train, y\_train, 'b', linewidth=3)

plt.plot(x\_test, y\_test, 'g', linewidth=3)

plt.legend(loc="upper left")

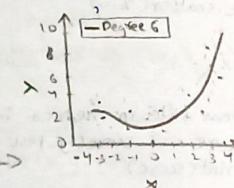
plt.xlabel('x')

plt.ylabel('y')

plt.axis([-4, 4, 0, 10])

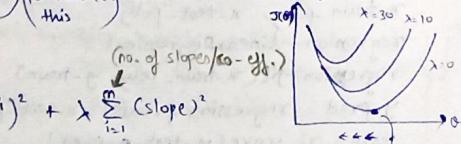
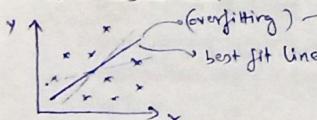
plt.show() # upper limit

poly\_regression(6)



## RIDGE, LASSO, & ELASTICNET ML ALGORITHMS

Ridge Regression: (or L2 regularization) → used to reduce overfitting



$$\text{Cost func.} = \frac{1}{2m} \sum_{i=1}^m (h_0(x)^i - y_i)^2 + \lambda \sum_{i=1}^n (\text{slope})^2$$

(this should never be equal to 0 as it would mean that it's overfitting)  $\propto \frac{1}{\text{slope}}$

$$\text{ex: } h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$(\text{let}) = 0.34 + 0.52x_1 + 0.48x_2 + 0.24x_3$$

| after Ridge

→ after applying ridge regression, value will decrease but it will never be equal to 0.

means; if  $x_1$  moves by 1, then 'y' will move by 0.52. Same for  $x_2$ .

Lasso Regression: (or L1 regularization) → used for feature selection

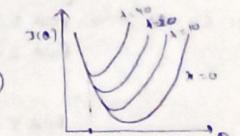
$$\text{cost func.} = \frac{1}{2m} \sum_{i=1}^m (h_0(x)^i - y_i)^2 + \lambda \sum_{i=1}^n |\text{slope}|$$

$$\text{ex: } h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 \quad (\text{4 independent features})$$

$$(122) = 0.52 + 0.65x_1 + 0.72x_2 + 0.39x_3 + 0.12x_4$$

| after lasso

$$= 0.52 + 0.51x_1 + 0.60x_2 + 0.14x_3 + 0.04x_4$$



$x_4$  is only 0.12 times related to output, so we can say that this feature is not that imp.

Here also 'θ' value decreases but after one point of time, it will become 0 resulting in removal of less-correlated features.

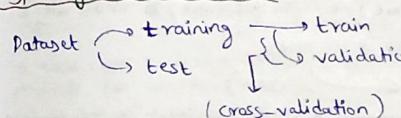
ElasticNet Regression:  $\begin{cases} \rightarrow \text{Reduce Overfitting} \\ \rightarrow \text{Feature Selection} \end{cases}$

$$\text{cost func.} = \frac{1}{2m} \sum_{i=1}^m (h_0(x)^i - y_i)^2 + \lambda_1 \sum_{i=1}^n (\text{slope})^2 + \lambda_2 \sum_{i=1}^n |\text{slope}|$$

reduce overfitting

feature selection

Types of Cross Validation:



1) Leave One Out CV:

ex: training → 500 records

exp1: Train → Acc 1

exp2: Train → Acc 2

⋮

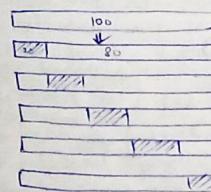
exp500: Train → Acc 500

Disadvantage: → complexity of training model increases

→ overfitting → train acc. ↑ ⚡ new test val. acc. ↓ data ↓ acc. ↓

2) K-Fold CV:

ex: K=5, n=100



Model trained on 5 different versions of training data & evaluated on 5 diff. versions of the test data.

avg. accuracies ⚡ test size =  $\frac{100}{5} = 20$

3) Stratified K-Fold CV:

It's similar to K-fold CV, but with an imp twist: It maintains the same proportion of different classes (labels) in each fold as in the entire dataset.

It's useful when working with imbalanced datasets.

#### 4) Time Series (V)

Here order of observations is crucial.

ex: cn = 4 → fold 1: train on [T1], test on [T2]

Fold 2: . . . [T1, T2], test on [T3]

, fold 5: . . . [T1, T2, T3, T4], test on [T5]

ex: financial forecasting, weather forecasting

Ex: (Ridge & Lasso Regression)

# importing

dataset = pd.read\_csv('Algerian\_forest\_fires\_dataset\_UPDATE.csv', header=1)

# .head(), .info()

# Data Cleaning

dataset.isnull().sum()

dataset[dataset.isnull().any(axis=1)]

dataset.loc[:, 'Region'] = 0 } creating new col.

dataset.loc[122:, 'Region'] = 1 }

df = dataset

df['Region'] = df['Region'].astype(int)

df.isnull().sum()

df = df.dropna().reset\_index(drop=True)

df = df.drop(122).reset\_index(drop=True) → refer (dataset) to know why?

df.columns → Note: There are whitespaces in col. names.

df.columns = df.columns.str.strip() → removing them

# change reg. col. to int dtype

df[['month', 'day', 'year', 'Temperature', 'RH', 'Ws']] = df[['month', 'day', 'year', 'Temperature', 'RH', 'Ws']].astype(int)

# change other col. to float dtype

objects = [features for features in df.columns if df[features].dtypes == 'U']  
for i in objects:

if i != 'classes':

df[i] = df[i].astype(float)

df.describe()

# save cleaned dataset

df.to\_csv('Aff-cleaned-dataset.csv', index=False)

# EDA

df\_copy = df.drop(['day', 'month', 'year'], axis=1)

df\_copy['classes'].value\_counts() → fire 131  
not fire 101 Additional spaces  
df\_copy['classes'] = np.where(df\_copy['classes'].str.contains('not fire'), 0, 1) fire 4 are present, so we have to remove that

# Visualizing

# Density Plot = Shows distribution of dataset over a continuous interval

plt.style.use('seaborn-v0\_8')

df\_copy.hist(bins=50, figsize=(20,15))

plt.show()

Correlation

df\_copy.corr()

sns.heatmap(df\_copy.corr(), numeric\_only=True)

# Pie chart

```
percentage = df['classes'].value_counts(normalize=True) * 100
classlabels = ['Fire', 'Not Fire']
plt.pie(percentage, labels=classlabels, autopct='%.1f%%')
plt.show()
```

# Box Plot

sns.boxplot(df['FWI'], color='green')

df['classes'] = np.where(df['classes'].str.contains('not fire'), 'not fire', 'fire')

# Monthly Fire Analysis

dftemp = df.loc[df['Region'] == 1]

sns.countplot(x='month', hue='classes', data=dftemp)

plt.title('Fire Analysis of Sidi-Bel-Regions', weight='bold')

- Similarly for df[Region] == 0

# Observations

Ex: (Model Training) → considering this as regression prob. → predicting 'FWI' col.

# Importing - reading dataset = 'Aff-cleaned-dataset.csv' - drop 'month', 'year', 'day'

# Encoding

df['classes'] = np.where(df['classes'].str.contains('not fire'), 0, 1)

df['classes'].value\_counts()

x = df.drop('FWI', axis=1) → Independent & Dependent Features

y = df['FWI'] → If we consider this as classification prob., we must remove 'classes' col.

②

x\_train.corr() → Note: Features which are highly (+)ve correlated can be removed

# check for multicollinearity

sns.heatmap(x\_train.corr(), annot=True)

# Func. to identify & return a set of col. names from the dataset that have corr. above a specific threshold

def correlation(dataset, threshold):

col\_corr = set()

corr\_matrix = dataset.corr()

for i in range(len(corr\_matrix.columns)):

for j in range(i+1):

if abs(corr\_matrix.iloc[i,j]) > threshold:

col\_name = corr\_matrix.columns[i]

col\_corr.add(col\_name)

return col\_corr

corr\_features = correlation(x\_train, 0.85)

x\_train.drop(corr\_features, axis=1, inplace=True)

x\_test.drop(corr\_features, axis=1, inplace=True)

x\_train.shape, x\_test.shape

# Standardization

① → x\_train\_scaled, x-test\_scaled

# Box plots to understand effect of StandardScaler

plt.subplots(figsize=(15,5))

plt.subplot(1,2,1)

sns.boxplot(data=x\_train)

plt.subplot(1,2,2)

sns.boxplot(data=x\_train\_scaled)

## # Linear Regression Model

```
from sklearn.linear_model import Lasso, Ridge, ElasticNet, LinearRegression  
from sklearn.metrics import mean_absolute_error, r2_score  
  
linreg = LinearRegression()  
linreg.fit(x_train_scaled, y_train)  
y_pred = linreg.predict(x_test_scaled)  
mae = mean_absolute_error(y_test, y_pred)  
score = r2_score(y_test, y_pred)  
plt.scatter(y_test, y_pred)
```

## # Lasso Regression

```
lasso = Lasso()  
⑤ → use 'lasso' instead of 'linreg'
```

$r^2$ -score is low doesn't mean the model is not performing well, we can also conclude that prev. model is overfitting.

## # Ridge Regression

```
ridge = Ridge()  
⑤ → use 'ridge' instead of 'linreg'
```

## # ElasticNet

```
elastic = ElasticNet()  
⑤ → use 'elastic' instead of 'linreg'
```

## # HyperParameter Tuning

```
from sklearn.linear_model import LassoCV, RidgeCV, ElasticNetCV  
lassocv = LassoCV(cv=10)  
lassocv.fit(x_train_scaled, y_train)  
lassocv.alpha_, lassocv.alpha_min_, lassocv.mse_path_  
⑤  
ridgecv = RidgeCV(cv=10)  
⑤  
elasticcv = ElasticNetCV(cv=10)  
⑤
```

(Explore all its features)

## STEP BY STEP PROJ. IMPLEMENTATION WITH LIFECYCLE OF ML PROJ.

(SEE NEXT PAGE)

Ex: Refer Linear Regression example (height-weight) in previous - previous section.

→ New data

```
scaled_weight = scaler.transform([[80]]) → scaled_weight = array([0.3135...])  
print("Height prediction for weight 80kg is: ", regression.predict(scaled_weight[0]))
```

# Assumptions

```
plt.scatter(y_test, y_pred_test) → If we are getting our points linearly,  
then we can assume we have done very good predictions
```

residuals = y-test - y-pred-test

```
sns.distplot(residuals, kde=True) → If residuals have a normal distribution,  
we can say we have created a good model.
```

```
plt.scatter(y_pred_test, residuals) → If we get uniform distribution, it confirms  
we have created a good model
```

Ex: (Multiple Linear Regression) ; Refer Notebook for op.

```
# Importing  
from sklearn.datasets import fetch_california_housing  
california = fetch_california_housing()  
california.keys(), print(california.DESCR), california.target_names  
california.feature_names, california.data.shape
```

# Prepare dataset

```
df = pd.DataFrame(california.data, columns=california.feature_names)  
df.head()
```

```
df['Price'] = california.target
```

```
= df.info(), .describe(), .corr(); sns.heatmap(df.corr(), annot=True)
```

# Independent & Dependent Features

```
x = df.iloc[:, :-1]
```

```
y = df.iloc[:, -1]
```

```
①
```

```
② → (slope/co-eff.)
```

```
regression.coef_, regression.intercept_
```

# Prediction for test data

```
y_pred = regression.predict(x_test)
```

```
②
```

# Assumptions

```
plt.scatter(y_test, y_pred)
```

```
residuals = y-test - y-pred
```

```
sns.distplot(residuals, kind='kde')
```

```
plt.scatter(y_pred, residuals)
```

Refer prev. page to know what they represent about model.

# Pickling : used for saving & loading a model.

```
import pickle  
pickle.dump(regression, open('regression.pkl', 'wb'))
```

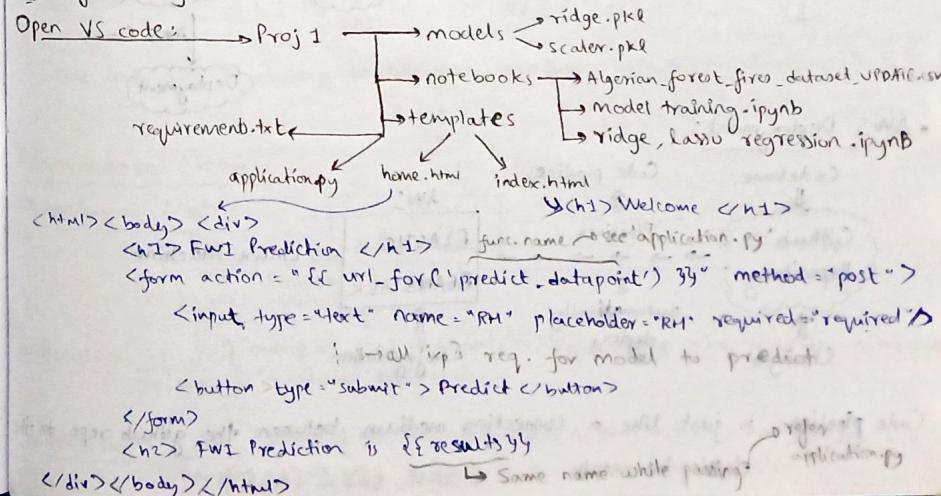
```
model = pickle.load(open('regression.pkl', 'rb'))
```

# Testing the saved model

```
model.predict(x_test)
```

→ Extension must be 'pkl'

Ex: Refer 'Ridge' Lasso & ElasticNet examples in prev. section



## application.py

```
# To install below libraries, run "pip install -r requirements.txt" in terminal
from flask import Flask, request, jsonify, render_template
import pickle

application = Flask(__name__) # name → "application" → so that it's easy to deploy it on AWS
app = application

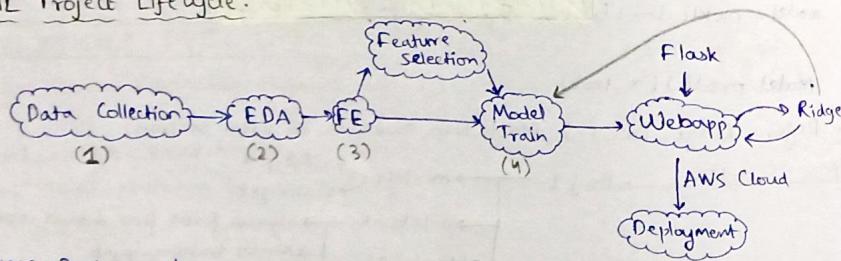
ridge_model = pickle.load(open('models/ridge.pkl', 'rb'))
standard_scaler = pickle.load(open('models/scaler.pkl', 'rb'))

@app.route('/')
def index():
    return render_template('index.html')

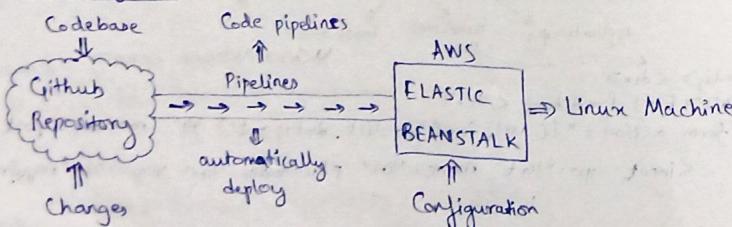
@app.route('/predictdata', methods=['GET', 'POST'])
def predict_datapoint():
    if request.method == 'POST':
        RH = float(request.form.get('RH'))
        all_ip = request.form.get('all_ip')
        new_data_scaled = standard_scaler.transform([[RH, all_ip]])
        result = ridge_model.predict(new_data_scaled) → result is in the form [0.0]
        return render_template('home.html', results=[result])
    else:
        return render_template('home.html')

if __name__ == '__main__':
    app.run(host='0.0.0.0')
# Deployed on AWS Elastic Beanstalk
```

## ML Project Lifecycle:



## AWS Deployment:

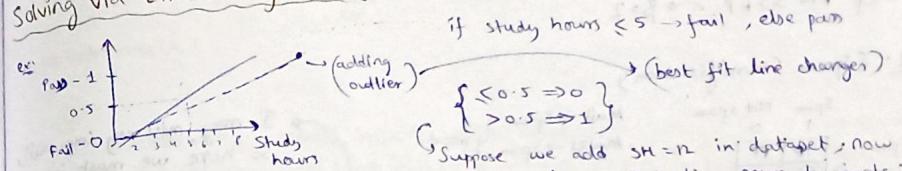


Code pipeline: is just like a connecting medium between the github repo. to the elastic beanstalk.  
 (SEE PREV. PAGE)

## LOGISTIC REGRESSION

Logistic Regression is used for solving "Binary Classification" problem.

### Solving via Linear Regression?

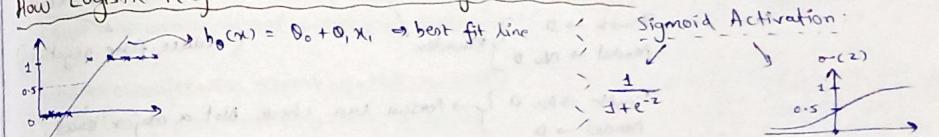


(Also, the val. is going <0 & >1, so we need to squash it i.e. once the line is out of bounds, make it parallel to X-axis)

Suppose we add SH=12 in dataset, now best fit line will incline more towards X-axis, & now SH=5 → ≤ 0.5, so it's not a good model

only possible with the help of Logistic Regression

### How Logistic Regression Solves Classification Prob.:



$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2 \quad \text{Linear Reg.} \rightarrow h_0(x) = \theta_0 + \theta_1 x \rightarrow \text{Convex func.}$$

let's denote:  $\text{cost}(h_0(x^{(i)}), y^{(i)})$

$$\text{log loss} \Rightarrow \text{cost}(h_0(x^{(i)}), y^{(i)}) = \begin{cases} -\log(h_0(x)) & \text{if } y=1 \\ -\log(1-h_0(x)) & \text{if } y=0 \end{cases} \quad \text{Linear Reg.} \rightarrow h_0(x) = \frac{1}{1+e^{-x}} \rightarrow \text{(non convex func.)}$$

↓ (proved mathematically)

$$\therefore \text{cost}(h_0(x^{(i)}), y^{(i)}) = -y^{(i)} \log(h_0(x^{(i)})) - (1-y^{(i)}) \log(1-h_0(x^{(i)}))$$

$$\therefore J(\theta_0, \theta_1) = -\frac{1}{2m} \sum_{i=1}^m (y^{(i)} \log(h_0(x^{(i)})) - (1-y^{(i)}) \log(1-h_0(x^{(i)})))$$

Our main aim is always to minimize cost func.  $J(\theta_0, \theta_1)$  by changing  $\theta_0$  &  $\theta_1$ .  
 (We can use convergence algo.)

### Performance Metrics:

#### 1) Confusion Matrix:

		Actual values	
		1	0
Predicted values	1	True(+)	False(+)
	0	False(-)	True(-)

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

2) Recall =  $\frac{TP}{TP+FN}$  } out of all predicted values, how many are correctly predicted  
 is shown by 'recall'.

3) Precision =  $\frac{TP}{TP+FP}$  } out of all correct results/actual values, how many are correctly predicted is shown by 'precision'.

$$4) F\text{-Beta Score} = \frac{(1+\beta^2)(\text{Precision} * \text{Recall})}{\text{Precision} + \text{Recall}}$$

- If FP & FN are both imp.  $\rightarrow \beta=1$
- If FP is more imp. than FN  $\rightarrow \beta=0.5$
- If FN >> FP  $\rightarrow \beta=2$

### Ex 1: Spam Classification

		Spam	Not spam	
Spam	1	TP	FP	Mail $\rightarrow$ Spam Model $\rightarrow$ Spam
	0	FN	TN	
		Mail $\rightarrow$ Not spam Model $\rightarrow$ Not spam		Blunder
		Mail $\rightarrow$ Not spam Model $\rightarrow$ Spam		

I may miss my imp. mail, so here we must increase our precision.  
(Reduce False Positives)

### Ex 2: Diabetes prediction

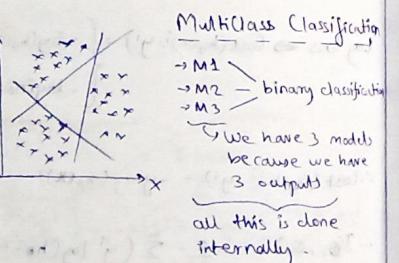
		D	No D	
D	1	TP	FP	Truth $\rightarrow$ D Model $\rightarrow$ D
	0	FN	TN	
		Truth $\rightarrow$ No D Model $\rightarrow$ No D		Blunder
		Truth $\rightarrow$ No D Model $\rightarrow$ D		person can check. Not a major issue!

Imp. for person  $\rightarrow$  so we must increase our recall  
(reduce false negatives)

### One vs One:

### One vs Rest (OVR):

			01	02	03	
01	01	1	0	0		MultiClass Classification
	02	0	1	0		
	03	0	0	1		
			01	1	0	



New test data  $\rightarrow$  M1  $\rightarrow$  0.25  
 $\rightarrow$  M2  $\rightarrow$  0.20  
 $\rightarrow$  M3  $\rightarrow$  0.55  $\rightarrow$   $0.55 = M3$

- OVR breaks down multi-class problem into several binary classification problems.
- In above ex., OVR creates 3 separate classifiers;

- class F1 vs Not F1 (i.e. F2 & F3)
- class F2 vs Not F2 (i.e. F1 & F3)
- class F3 vs Not F3 (i.e. F1 & F2)

```
# importing
from sklearn.datasets import make_classification
# create dataset
X, y, = make_classification(n_samples=1000, n_features=10,
n_classes=2, random_state=15)
# Creates standardized data
# so we need not do it.
```

pd.DataFrame(X)

y

①

```
from sklearn.linear_model import LogisticRegression
logistic = LogisticRegression() # Shift + tab few times to see these as well.
logistic.fit(x_train, y_train)
y_pred = logistic.predict(x_test)
print(y_pred)
logistic.predict_proba(x_test) # probability of outputs
```

```
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
score = accuracy_score(y_test, y_pred)
```

```
print(classification_report(y_test, y_pred))
print(confusion_matrix(y_test, y_pred))
```

### # Hyperparameter Tuning

```
model = LogisticRegression()
params = {'penalty': ['l1', 'l2', ...],
'C': [100, ...],
'solver': ['sag', ...]}
```

```
from sklearn.model_selection import GridSearchCV, StratifiedKFold, RandomizedSearchCV
cv = StratifiedKFold() # we can use K-fold cv as well.
```

```
grid = GridSearchCV(estimator=model, param_grid=params, scoring='accuracy', cv=cv,
grid.fit(x_train, y_train)
grid.best_params_ # grid.best_score_
y_pred = grid.predict(x_test)
```

### # RandomSearch

```
⑨ → change : ⑧ grid → randomcv ⑩ GridSearchCV → RandomizedSearchCV
```

### # For multiclass classification problem

```
x, y = make_classification(n_samples=1000, n_features=10, n_informative=3,
n_classes=3, random_state=1)
```

```
⑦ → change to: logistic = LogisticRegression(multi_class='ovr')
```

⑥

### # For Imbalanced dataset

```
from collections import Counter
# by def. 2 or features
x, y = make_classification(n_samples=1000, n_features=2, n_clusters_per_class=1,
n_redundant=0, weights=[0.99], random_state=10)
```

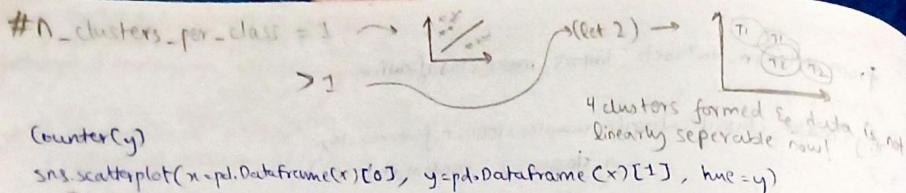
No extra redundant

features are added.

All features are informative, no duplicates.

→ controls how many groups are formed within each class.

99% of samples belong to class 0, making it imbalanced



(Counter(y))

sns.scatterplot(x=pd.DataFrame(x)[0], y=pd.DataFrame(x)[1], hue=y)

⑤

⑥ → added few more parameters

⑦

⑧

# With ROC Curve and ROC AUC Score → Plot of TPR vs FPR at various threshold levels, showing a model's performance across different threshold values.

# Now "model.predict(y-test)" gives o/p in 0s & 1s. By definition, if threshold val. = 0.5 means o/p of 1 → th>0.5. Now depending on different use cases, we need to change th<0.5

from sklearn.metrics import roc\_curve, roc\_auc\_score

from matplotlib import pyplot

x,y = make\_classification(n\_samples=1000, n\_classes=2, random\_state=42)

⑨

dummy\_model\_prob = [0 for \_ in range(len(y-test))] → dummy model with default output as 0

Model = LogisticRegression()

Model.fit(x-train, y-train)

Model\_prob = Model.predict\_proba(x-test) → gives probabilities of output feature

Model\_prob

Model\_prob = Model\_prob[:, 1] → prob. for class 1 ← we are selecting the 2nd column

dummy\_model\_auc = roc\_auc\_score(y-test, dummy\_model\_prob)

Model\_auc = roc\_auc\_score(y-test, Model\_prob)

print() → 0.5 → obvio. → since dummy, so half of them must be true

0.907648

dummy\_fpr, dummy\_tpr, \_ = roc\_curve(y-test, dummy\_model\_prob)

Model\_fpr, Model\_tpr, thresholds = roc\_curve(y-test, Model\_prob)

Model\_fpr, Model\_tpr, thresholds → returns FPR, TPR, Threshold

(for ignoring 'size')

Pyplot.plot(dummy\_fpr, dummy\_tpr, linestyle='--', label='Dummy Model')

Pyplot.plot(Model\_fpr, Model\_tpr, marker='.', label='Logistic')

Pyplot.xlabel('FPR')

Pyplot.ylabel('TPR')

Pyplot.legend()

Pyplot.show()

After we plot this curve, we have to select a threshold which results in high TPR & low FPR. Once that point is detected, we will select that threshold for our Logistic Regression

Display threshold values

fig = Pyplot.figure(figsize=(20, 50))

ax = fig.subplot(111)

for xyz in zip(Model\_fpr, Model\_tpr, thresholds):

ax.annotate(f'({xyz[0]}, {xyz[1]}), xyz[2])

sample	actual class	predicted prob.	classification (thresh 0.5 = 0.5)
1	1	0.95	+
2	1	0.85	+
3	1	0.6	-
4	1	0.45	-
5	0	0.8	+
6	0	0.9	+
7	0	0.3	-

See above table for better understanding of ROC curve.

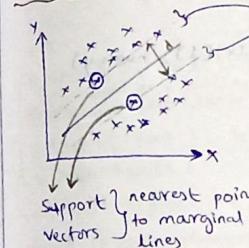
Also, threshold value is told by "domain expert".

## SUPPORT VECTOR MACHINES

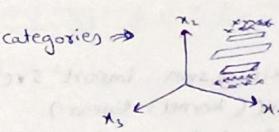
SVM is ML Algo. helps us solve both classification & regression problems.

### 1) Support Vector Classifier (SVC):

marginal planes → created in such a way that distance (d) btw them is maximum



if we had 3 categories →



Hard Margin: requires a strict separation with no misclassifications. suitable for linearly separable data.

### Soft Margin:

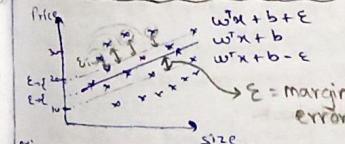
Allows some misclassification, making it more flexible for non-linearly separable data. In reality, there will always be misclassification.

$$\text{Cost func.} = \min_{(w, b)} \frac{\|w\|}{2} + C \sum_{i=1}^n \xi_i \quad \text{(Spelled as 'lata')}$$

(Soft Margin)

(how many points we want to avoid misclassification) → summation of the distance of the incorrect data points from the marginal plane.

### 2) Support Vector Regression (SVR):



$$\text{Cost func.} = \min_{(w, b)} \frac{\|w\|}{2} + C \sum_{i=1}^n \xi_i \quad \text{(Spelled as 'lata')}$$

constraint:  $|y_i - w^T x_i| \leq \epsilon + \xi_i$

loss func. (margin) error (above the margin)

Note:  $C \propto \frac{1}{\text{Loss func.}}$

↓ Loss func.

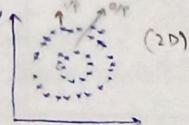
↓ C

↓ Loss func.

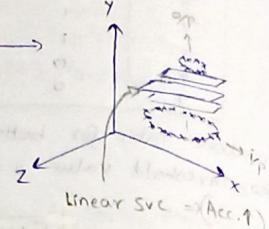
↓ C

## SVM Kernels

When datapoints are not linearly separable, we use SVM kernels



→ Transformations  
(We apply some mathematical formulas)



We have different types of kernels:

- Polynomial Kernel
- RBF Kernel
- Sigmoid Kernel

Ex: (SVC Implementation)

# importing

```
from sklearn import svm
Xy = make_classification(n_samples=1000, n_features=2, n_classes=2,
                        n_clusters_per_class=2, n_redundant=0)
```

x, y

```
sns.scatterplot(x=pd.DataFrame(x)[0], y=pd.DataFrame(x)[1], hue=y)
# Points at 'x' & 'y' → point color based on class of 'hue' → class 0 → Blue Colour
# class 1 → Orange Colour
```

④

```
from sklearn.svm import SVC
svc = SVC(kernel='linear')
svc.fit(x_train, y_train)
svc.coef_
y_pred = svc.predict(x_test)
```

⑤

# Similarly trying for other kernels

→ rbf = SVC(kernel='rbf')

→ polynomial = SVC(kernel='poly')

→ sigmoid = SVC(kernel='sigmoid')

# Hyperparameter Tuning

```
param_grid = {'C': [0.1, 1, 10],
              'gamma': [1, -1],
              'kernel': ['rbf', -3]}
```

- controls the amount of info. to be displayed in console during execution of grid search process.  
- Higher values = more detailed logs

```
grid = GridSearchCV(SVC()), param_grid = param_grid,
grid.fit(x_train, y_train)
grid.best_params_
y_pred2 = grid.predict(x_test)
print(classification_report(y_test, y_pred2))
print(confusion_matrix(y_test, y_pred2))
```

means after finding the best combo. of hyperparams based on cross-validation, the model is retrained (refit) on the entire training dataset using those best hyperparams.

Ex: (SVM Kernels) In-depth Explanation → Refer Notebook

Basically → initial dataset ⇒ ↑(0)

→ Creating few new features & visualizing using "scatter\_3d"

→ Using different kernels of SVC → fit

→ predict

→ performance metrics

(SVR Implementation)

```
Ex: (SVR Implementation)
# importing
df = sns.load_dataset('tips')
df.info(); df['ccol.name'].value_counts() → checking for diff. columns
df.columns
df.independent, df.dependent_features
x = df[['tip', 'sex', 'smoker', 'day', 'time', 'size']]
y = df['total_bill']
```

⑥

# Feature Encoding

(here 'days'  
col. has  
4 values)

Label encoding: If we have 2 values, we can substitute by 0 and 1.

OHE: If > 2 values, we can create those many features & wherever that specific value is, it will be = 1 & rest all = 0.

```
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
```

le1 = LabelEncoder() } applying Label Encoding to binary features

le2 = "

le3 = "

x\_train['sex'] = le1.fit\_transform(x\_train['sex'])

x\_train['smoker'] = le2.fit\_transform(x\_train['smoker'])

x\_train['Time'] = le3.fit\_transform(x\_train['Time'])

x\_train.head()

→ (...) replace 'x\_train' with 'x\_test'

```
from sklearn.compose import ColumnTransformer → OHE applied with the help of this
```

ct = ColumnTransformer(transformers=[('onehot', OneHotEncoder(drop='first'), [3])], remainder='passthrough')

(takes value in the form of tuples)

ex: RGB. We drop 'R' col which is first col. & now if both G=0 & B=0, we can say that it represents 'R' col.

'days' = 3rd index in table (starting from 0)

(We do this for efficiency & it also avoids multicollinearity & avoids redundant columns without losing any info.)

import sys

import numpy as np

np.set\_printoptions(threshold=sys.maxsize) → Just a setting to see the entire dataset

ct.fit\_transform(x\_train)

x\_train = ct.fit\_transform(x\_train)

x\_test = ct.transform(x\_test)

```
from sklearn.svm import SVR
```

svr = SVR()

svr.fit(x\_train, y\_train)

y\_pred = svr.predict(x\_test)

⑦

# Hyperparameter Tuning

⑧ → change 'SVC' to 'SVR()'

print(r2\_score(y\_test, y\_pred2))

print(mean\_absolute\_error(y\_test, y\_pred2))

## NAIVE BAYES THEOREM

Naive Bayes is used for "Classification" problem.

Bayes' Theorem:

$$\{P(A \text{ and } B) = P(A) * P(B|A)\}$$

$$P(A/B) = \frac{P(A) * P(B/A)}{P(B)}$$

$P(A/B)$  = prob. of event A, given B has occurred

$P(A)$  = prob. of event A,  $P(B)$  = prob. of event B

$P(B/A)$  = prob. of event B, given A has occurred

Variants of Naive Bayes:

1) Bernoulli Naive Bayes: Whenever your features are following a Bernoulli distribution, then we need to use Bernoulli Naive Bayes Algorithm.

Refer "Mathematics - Basics to Advanced for Data Science & GenAI" chapter

2) Multinomial Naive Bayes: when 'ip' is in text format

↓ converted to numerical values ↓ with the help of NLP  
 → ① Bow  
 → ② Tf-IDf  
 → ③ word2vec  
 ↓ vectors

3) Gaussian Naive Bayes: If the features are following gaussian distribution i.e. we have continuous values, then we use this alg.  
 ↗ 'IRIS' dataset ~ is a continuous dataset

Ex:

```
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
x, y = load_iris(return_X_y=True)
# Capitalized
```

```
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=42)
from sklearn.naive_bayes import GaussianNB
```

```
gnb = GaussianNB()
gnb.fit(x_train, y_train)
```

```
y_pred = gnb.predict(x_test)
```

```
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
print(confusion_matrix(y_pred, y_test))
print(accuracy_score(y_pred, y_test))
print(classification_report(y_pred, y_test))
```

// Refer how to do feature engineering in a classification problem.

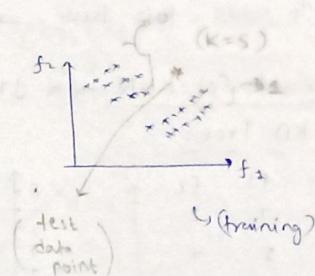
## K NEAREST NEIGHBOUR ML ALGORITHM

KNN ML Algo. helps us solve both classification & regression problems.

Classification:

Steps: 1) Initialize 'K' value; ( $K > 0$ )  
 $K = 1, 2, 3, \dots \Rightarrow$  hyperparameter

- 2) find the K Nearest Neighbours for the test data
- 3) from those, ( $K=5$ ) how many neighbour belong to 0 & 1 category.



We find it's 5 nearest neighbours & conclude on basis of that.

2 ways to calculate distance:

1) Euclidean Distance:

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

2) Manhattan Distance:

$$d = |x_2 - x_1| + |y_2 - y_1|$$

Ex: # importing

```
my = make_classification(n_samples=1000, n_features=3, n_redundant=1,
n_classes=2, random_state=42)
```

①

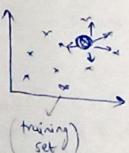
```
from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier(n_neighbors=5, algorithm='auto')
classifier.fit(x_train, y_train)
```

```
y_pred = classifier.predict(x_test)
```

⑥

{ Try Hyperparameter Tuning using GridSearchCV for diff. 'K' values }

2) Regression:



For test data, we calc. → avg. of all these 'K' nearest points to find the o/p.

If we have a lot of outliers, we can go with median.

Ex: # importing

```
from sklearn.datasets import make_regression
x, y = make_regression(n_samples=1000, n_features=2, noise=10,
random_state=42)
```

①

```
from sklearn.neighbors import KNeighborsRegressor
```

```
regressor = KNeighborsRegressor(n_neighbors=6, algorithm='auto')
```

```
regressor.fit(x_train, y_train)
```

```
y_pred = regressor.predict(x_test)
```

②

## Issue



To find distance

$$T_c = O(n^2) \rightarrow \text{if } n \sim \text{million} \rightarrow \text{High time complexity}$$

To solve this issue, we have 2 optimization

KD Tree  
Ball Tree

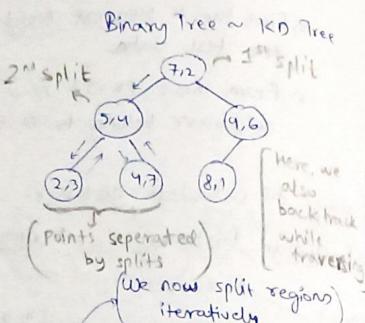
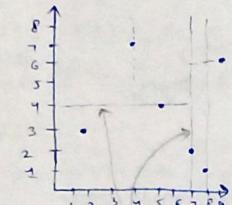
## Variants of KNN $\rightarrow$ (to $J T_c$ )

### 1) KD Tree:

$f_1$	$f_2$
7	2
5	4
9	6
2	3
4	7
8	1

$$\text{Median} = 5.5 \quad 3.5$$

(upper value which is there in dataset)

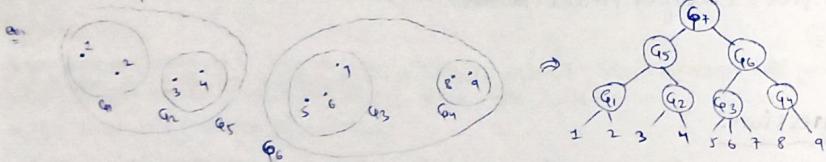


Now we split regions  $\rightarrow$  this reduces search area for nearest neighbors.

Points left of median go into one region; points to the right go into another

### 2) Ball Tree:

Recursively divide dataset into "balls", where each ball consists a subset of the data points.



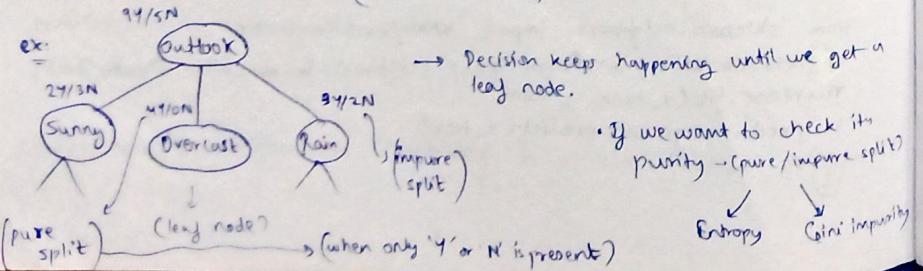
## DECISION TREE CLASSIFIER & REGRESSOR

### 1) Classification:

Decision Tree Classifier

$\rightarrow$  ID3 (Iterative Dichotomiser) makes 'generic tree'

$\rightarrow$  CART (sklearn uses CART) makes 'binary tree'

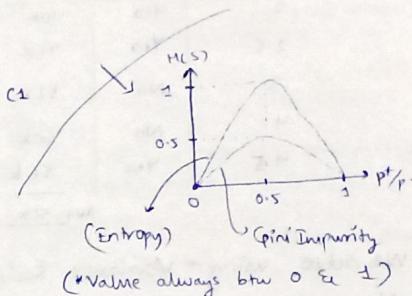


what feature we need to select for splitting  $\rightarrow$  Info. gain

$$\text{Entropy: } H(S) = -P_1 \log_2(P_1) - P_2 \log_2(P_2)$$

$$H(C_1) = -\frac{3}{6} \log_2(\frac{3}{6}) - \frac{3}{6} \log_2(\frac{3}{6}) = 1 \rightarrow \text{Impure split}$$

$$H(C_2) = -\frac{3}{3} \log_2(\frac{3}{3}) - 0 \cdot \log_2(0) = -1 \log_2 1 = 0 \rightarrow \text{Pure split}$$



$$\text{Gini Impurity: } G.I. = 1 - \sum_{i=1}^k (P_i)^2$$

$$= 1 - (C_{P_1})^2 + (C_{P_2})^2 \rightarrow \text{For } C_1$$

$$G.I. = 1 - (\frac{3}{3})^2 = 1 - (1/2)^2 + (1/2)^2$$

$$= 0 \rightarrow \text{Pure split}$$

$$\text{If 3 categories} \rightarrow \text{Gini Impurity}$$

$$H(S) = -P_1 \log_2(P_1) - P_2 \log_2(P_2) - P_3 \log_2(P_3)$$

Whenever dataset is small  $\rightarrow$  use "Entropy" (determine quality of split)  
" " " " large  $\rightarrow$  use "Gini Impurity" (default by decision tree)

Information Gain: (helps us understand which feature one should select to do the decision tree split)

$$\text{Gain}(S, f_1) = H(S) - \sum_{\text{value}} \frac{|S_v|}{|S|} H(S_v)$$

$$\rightarrow H(S) = -P_1 \log_2(P_1) - P_2 \log_2(P_2) \\ = -\frac{3}{6} \log_2(\frac{3}{6}) - \frac{3}{6} \log_2(\frac{3}{6}) \approx 0.94$$

$$H(C_1) = -\frac{3}{6} \log_2(\frac{3}{6}) - \frac{3}{6} \log_2(\frac{3}{6}) \approx 0.81$$

$$H(C_2) = 1$$

$$\therefore \text{Gain}(S, f_1) = 0.94 - [(\frac{3}{6} \times 0.81) + (\frac{3}{6} \times 1)] = 0.649$$

$\therefore \text{Gain}(S, f_1) = 0.649 > \text{Gain}(S, f_2) = 0.51$

This indicates that I have to start splitting from feature f2.

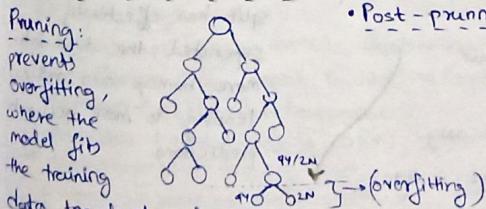
### Post & Pre Pruning:

#### Pruning:

Prevents overfitting, where the model fits the training data too closely, leading to 'high variance' in test data.

• Post-pruning: Builds complete decision tree first, then prune unnecessary branches.

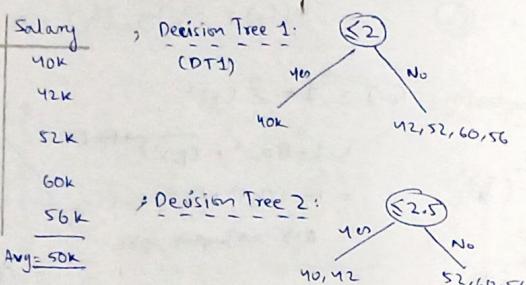
Effectively reduces overfitting  
Best for smaller datasets



- Pre-pruning:
  - limits the tree's growth during construction by setting hyperparameters like 'max\_depth', etc.
  - Ideal for large datasets

## 2) Regression:

Ex 1	Experience	gap	Salary
2	Yes	400	40K
2.5	Yes	420	42K
3	No		52K
4	No		60K
4.5	Yes	500	56K



We solve using "Variance Reduction" =  $\text{Var}(\text{parent}) - \sum w_i \cdot \text{var}(\text{child}_i)$

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad (\text{MSE})$$

$$\therefore \text{Var}(\text{root}) = \frac{1}{5} [(40-50)^2 + (42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2] = 60.8$$

$$\text{For DT1: } \therefore \text{Var}(C1) = \frac{1}{2} \sum_{i=1}^2 (y_i - \bar{y})^2 = \frac{1}{2} (40-50)^2 = 100$$

$$\therefore \text{Var}(C2) = \frac{1}{3} \sum_{i=1}^3 (y_i - \bar{y})^2 = \frac{1}{3} [(42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2] = 51$$

$$\therefore \text{Variance Reduction} = \text{Var}(\text{root}) - \sum w_i \cdot \text{Var}(\text{child}_i) = 60.8 - \left[ \frac{1}{2} \times 100 + \frac{1}{3} \times 51 \right] = 60.8 - 20 - 40.8 = 0$$

$$\therefore (VR)_1 = 0$$

$$\text{For DT2: } \therefore \text{Var}(C1) = \frac{1}{2} [(40-50)^2 + (42-50)^2] = 82$$

$$\therefore \text{Var}(C2) = \frac{1}{3} [4 + 100 + 36] = 46.66$$

$$\therefore (VR)_2 = 60.8 - \left( \frac{2}{5} \times 82 + \frac{3}{5} \times 46.66 \right) = +0.004$$

$\therefore (VR)_1 < (VR)_2 \rightarrow \therefore$  We select 2nd split → as it indicates that the

split has effectively separated the data into more homogenous groups, leading to more accurate predictions.

So now → if test data →  $\leq 2.5 = \text{true}$

↓

we will take avg.

$$\text{i.e. } \text{avg} = \frac{40+42}{2} = 41$$

## Ex: (Decision Tree Classifier)

```
# importing
from sklearn.datasets import load_iris
iris = load_iris()
print(iris['DESCR'])
print(iris['target'])

# independent features
X = pd.DataFrame(iris['data'], columns=['sepal length', 'sepal width', 'petal length', 'petal width'])
y = iris['target']
```

① from sklearn.tree import DecisionTreeClassifier  
 tree\_classifier = DecisionTreeClassifier()  
 tree\_classifier.fit(X\_train, y\_train) → Post-pruning. we can set diff. parameters here. For ex: 'max\_depth=2'  
 # visualize (not affected by scale of features so no need to apply standardization)

from sklearn import tree  
 plt.figure(figsize=(15,10)) → ensures nodes in decision tree are coloured  
 tree.plot\_tree(tree\_classifier, filled=True) → because we have 3 distinct classes of iris flowers: Setosa, Versicolour & Virginica

② → Conf. matrix =  $3 \times 3 = \begin{bmatrix} TP_1 & FP_1 & FP_2 \\ FN_1 & TP_2 & FP_3 \\ FN_2 & FN_3 & TP_3 \end{bmatrix}$   
 # Prepruning & Hyperparameter Tuning

param = { 'criterion': ['gini', ...], 'splitter': ['best', ...], 'max\_depth': [1, 2, ...], 'max\_features': ['auto', ...] }

from sklearn.model\_selection import GridSearchCV  
 grid = GridSearchCV(tree\_classifier, param\_grid=param, cv=5, scoring='accuracy')  
 import warnings  
 warnings.filterwarnings('ignore')  
 grid.fit(X\_train, y\_train)  
 grid.best\_params\_, grid.best\_score\_

③ y\_pred = grid.predict(X\_test)

## Ex: (Decision Tree Regressor → Diabetes Prediction)

```
from sklearn.datasets import load_diabetes
dataset = load_diabetes()
print(dataset['DESCR'])
dataset

df_diabetes = pd.DataFrame(dataset.data, columns=['age', 'sex', 'bmi', 'bp', 's1', 's2', 's3', 's4', 's5', 's6'])

df_diabetes.head()
X = df_diabetes, y = dataset['target']
```

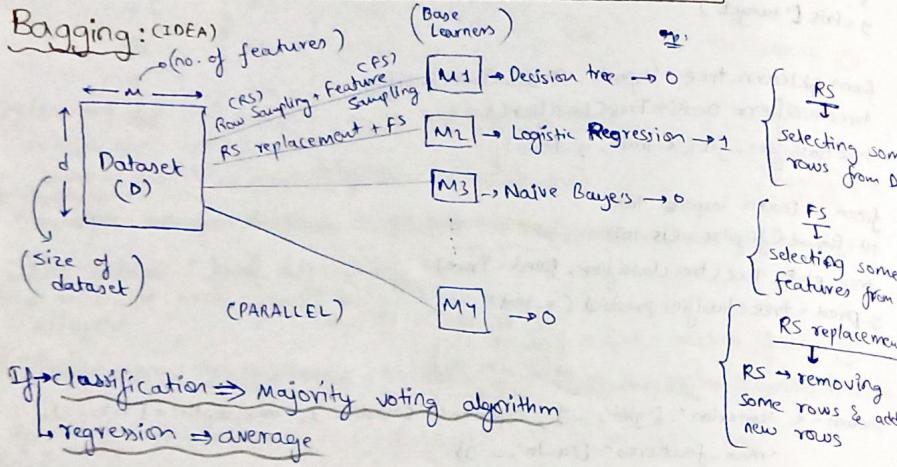
④ X\_train, corr()  
 plt.figure(figsize=(15,10))  
 sns.heatmap(X\_train.corr(), annot=True)  
 from sklearn.tree import DecisionTreeRegressor  
 regressor = DecisionTreeRegressor()  
 regressor.fit(X\_train, y\_train)  
 # Hyperparameter tuning  
 param → same as above one  
 grid = GridSearchCV(regressor, param\_grid=param, cv=5, scoring='neg\_mean\_squared\_error')

## // Plotting tree

```
selected model = DecisionTreeRegressor(criterion='squared_error', max_depth=4, max_features=1, splitter='random')
selected model .fit(x_train, y_train)
from sklearn import tree
plt.figure(figsize=(15, 15))
tree.plot_tree(selected model, filled=True)
```

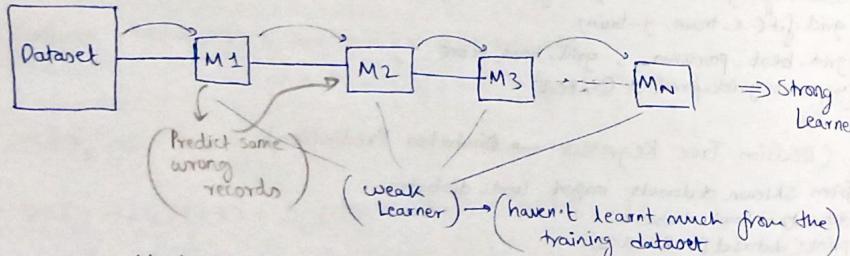
got from 'grid.best\_params\_'

## RANDOM FOREST MACHINE LEARNING



## AdaBoost

Boosting: (IDEA)



Process:

- M1 trains on dataset. It predicts correct & wrong gp's.
- Now M1 passes these wrong records along with some training dataset to M2.
- Above process repeats. (SEQUENTIALLY)

\* Now in Random forest → we use all Models (M1, M2, ..., MN) = Decision Trees

Q) Why should we use Random forest instead of DT?

### Decision Tree

↓  
Overfitting

train acc. ↑ = Low bias → Low bias  
test acc. ↓ = High variance → Low variance

### Random Forest

Generalized model

Ex: (Random Forest Classification)

# Importing

df = pd.read\_csv('Train.csv')

df.head()

# Feature Engineering → Data Cleaning → Handling Missing Values

df.isnull().sum()

df['Gender'].value\_counts() → similarly check for other categories

features\_with\_na = [feature for feature in df.columns if df[feature].isnull().sum() >= 1]

for feature in features\_with\_na:  
print(feature, np.round(df[feature].isnull().mean() \* 100, 2), '% missing values!')

# Statistics on numerical col

df[features\_with\_na].select\_dtypes(exclude='object').describe()

# Inputting null values

df.Age.fillna(df.Age.median(), inplace=True)

df.TypeofContract.fillna(df.TypeofContract.mode()[0])  
↳ (col.name)

df.isnull().sum()  
(used to access the 1st/most frequent val.)

df.drop('CustomerID', inplace=True, axis=1) → Not an input prop., so removing it

# Feature Extraction

df['TotalVisiting'] = df['NumberOfPersonVisiting'] + df['NumberOfChildrenVisiting']

df.drop(columns=['( )', '( )'], axis=1, inplace=True)

df.shape

# Get all numeric features

num\_features = [feature for feature in df.columns if df[feature].dtype != 'O']  
print(num\_features)

# Categorical features

cat\_features = [feature for feature in df.columns if df[feature].dtype == 'O']  
print(cat\_features)

# Discrete features

discrete\_features = [feature for feature in num\_features if len(df[feature].unique()) < 2]  
print(discrete\_features)

# Continuous features

continuous\_features = [feature for feature in num\_features if feature not in discrete\_features]  
print(continuous\_features)

# Independent & dependent features

x = df.drop(['ProdTaken'], axis=1)

y = df['ProdTaken']

x.head(), y.value\_counts()

①

x.info()

cat\_features = x.select\_dtypes(include='object').columns

num\_features = x.select\_dtypes(include='object').columns

from sklearn.preprocessing import OneHotEncoder, StandardScaler

from sklearn.compose import ColumnTransformer

numeric\_transformer = StandardScaler()

oh\_transformer = OneHotEncoder(drop='first')

preprocessor = ColumnTransformer([

(‘OneHotEncoder’, oh\_transformer, cat\_features)]

(‘StandardScaler’, numeric\_transformer, num\_features)])

Preprocessor

[Refer "Support Vector Machines"]



```

numeric_transformer = StandardScaler()
oh_transformer = OneHotEncoder(drop='first')
preprocessor = ColumnTransformer([
    ("OneHotEncoder", oh_transformer, onehot_columns),
    ("StandardScaler", numeric_transformer, num_features),
], remainder='passthrough')

```

X = preprocessor.fit\_transform(x)  
pd.DataFrame(x)

## # Model Training & Selection

# Importing: RandomForestRegressor, KNeighborsRegressor, Ridge, Lasso, DecisionTreeRegressor,  
sklearn.ensemble    sklearn.neighbors    sklearn.linear\_model    sklearn.tree

- (11) → initialize above models  
→ for performance → use: "RMSE", "MAE", "R2 score"  
→ rest all same.

## # Hyperparameter Tuning

knn\_params = {"n\_neighbors": [2, 3, 10, 20, 40, 50]}

rf\_params → Same as prev. example

randomcv\_models = [{"KNN": KNeighborsRegressor(), knn\_params},  
{"RF": RandomForestRegressor(), rf\_params}])

- (12) → add KNN in models' as well. Use best\_params\_ from above cell if p.  
(13) → for performance → use: "RMSE", "MAE", "R2 score"  
→ rest all same

## ADABOOST MACHINE LEARNING ALGORITHM

### Ensemble Techniques

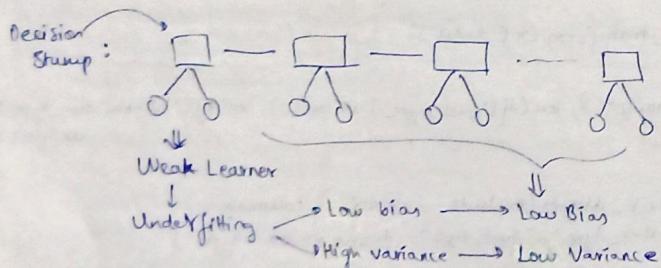
Bagging  
→ Random Forest Classifier  
→ Random Forest Regressor

Boosting  
→ Adaboost  
→ Gradient Boosting  
→ Xgboost

→ Adaboost: We assign "weights" to weak learners.

$$\text{func. } f = \alpha_1(M_1) + \alpha_2(M_2) + \dots + \alpha_n(M_n); \quad \{\alpha_1, \alpha_2, \dots, \alpha_n\}: \text{weights}$$

$\{M_1, M_2, \dots, M_n\}: \text{Decision Tree Stumps}$



### Maths. Intuition:

	Credit	Approval	Samplewt
≤ 50K	B	No	$\frac{1}{2}$
<= 50K	G	Yes	$\frac{1}{2}$
<= 50K	G	Yes	$\frac{1}{2}$
<= 50K	B	No	$\frac{1}{2}$
> 50K	G	Yes	$\frac{1}{2}$
> 50K	N	Yes	$\frac{1}{2}$
> 50K	N	No	$\frac{1}{2}$

② Sum of total error & performance of Stump? ←

$$(2.1) \text{ Sum of all the total error} = \frac{1}{2} = (\text{TE})$$

$$(2.2) \text{ Performance of Stump} = \frac{1}{2} \ln \left[ \frac{1 - \text{TE}}{\text{TE}} \right]$$

$$= \frac{1}{2} \ln \left( \frac{1 - \frac{1}{2}}{\frac{1}{2}} \right) \approx 0.896$$

$$\therefore f = \alpha_1(M_1) + \alpha_2(M_2) + \dots + \alpha_n(M_n)$$

$$\downarrow \\ (0.896) \rightarrow \text{Weight}$$

③ Update the weights for correctly & incorrectly classified points

↳ whichever is correctly classified, we reduce their weights & whichever is incorrectly classified, we increase their weights. We do this because there will be a high probability of selecting the wrong records by the next decision tree stump.

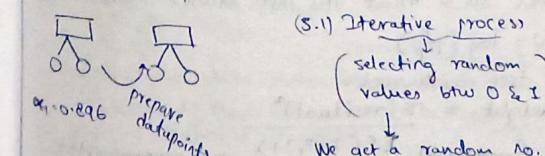
### Normalized Weights Computation

Salary	Credit	Approval	Update Wts	Normalized Wts
≤ 50K	B	N	0.058	0.08
<= 50K	G	Y	0.058	0.08
<= 50K	G	Y	0.058	0.08
> 50K	B	N	0.058	0.08
> 50K	G	Y	0.058	0.08
→ > 50K	N	Y	0.349	0.50
> 50K	N	N	0.058	0.08

$\sum = 0.697$        $\sum = 1$

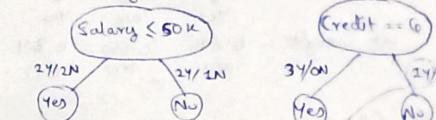
(But sum must be 1, so we normalize it) → (divide by 0.697)

④ Select data points to send to next stump



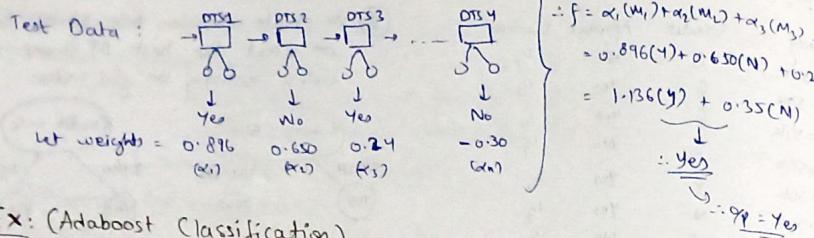
We get a random no., save check in which 'bin' it lies. Note: Rows can be repeated.

① We create decision tree stump & we select the best DT stump using entropy or gini index.



When credit = G, we have all 'Y'. Here we are getting 'Y' when credit = G, we can say this is an error.

## ④ Final Prediction



## Ex: (Adaboost Classification)

Refer "Random Forest Classification" example. Keep everything same.  
Just add "from sklearn.ensemble import AdaBoostClassifier" while training & checkout the results. You may also hyperparameter tune it. {Refer docs}

## Ex: (AdaBoost Regression)

Refer "Random Forest Regressor" example. Keep everything same.  
Just add "from sklearn.ensemble import AdaBoostRegressor" while training & checkout the results. You may also hyperparameter tune it. {Refer docs}

## GRADIENT BOOSTING

Steps: 1) Create a base model

2) Compute residuals, error  $\rightarrow r_i$  (Ex: Compute avg. & then subtract  $y_i$ )

3) Construct a decision tree considering  $y_p = x_i$  &  $y_p = r_i$

$$\therefore F(x) = \alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots + \alpha_n h_n(x) \quad ; \alpha_i \rightarrow \text{learning rate}$$

$$\therefore F(x) = \sum_{i=0}^n \alpha_i h_i(x) \quad \rightarrow \text{Final func. of gradient boosting}$$

## Ex: (GradientBoost Classification)

Refer "Random Forest Classification" example. Keep everything same.  
Just add "from sklearn.ensemble import GradientBoostingClassifier" while training & checkout the results. You may also hyperparameter tune it. {Refer docs}

## Ex: (GradientBoost Regressor)

Refer "Random Forest Regressor" example. keep everything same.  
Just add "from sklearn.ensemble import GradientBoostingRegressor" while training & checkout the results. You may also hyperparameter tune it. {Refer docs}

## XGBOOST MACHINE LEARNING ALGORITHM

### 1) Classification:

Steps: 1) Construct a base model such that it's not biased towards anything.  $\{ \log(\text{odds}) = \log \left( \frac{P}{1-P} \right) \}$

2) Construct a decision tree with root.

3) Calculate similarity weight =  $\sum (\text{residual})^2$

(Cover value)  $\sum P_i(1-P_i) + \lambda$   $\rightarrow$  Hyperparameter

4) Calculate Gain

5) Wrt finding final esp, we apply sigmoid activation func.

Ex: Refer "Random Forest Classification" example. Keep everything same.  
Just add "from xgboost import XGBClassifier" while training & checkout the results.  $\{ \text{if not present } \rightarrow !\text{pip install xgboost} \}$   
You may also hyperparameter tune it. {Refer docs}

## 2) Regressor:

Create a base model.

Steps: 1) Residual Computation.

2) Construct Decision Tree 1 using  $\{x_i, R_i\}$

Ex: Refer "Random Forest Regressor" example. Keep everything same.

Just add "from xgboost import XGBRegressor" while training & checkout the results. You may also hyperparameter tune it. {Refer docs}

## UNSUPERVISED MACHINE LEARNING

Unsupervised ML  $\rightarrow$  Clustering: group your data into similar clusters.

1) K Means Algo. 2) Hierarchical Clustering 3) DBScan Clustering 4) Silhouette Clustering

In supervised ML  $\rightarrow$  we have esp features / dependent features, but in unsupervised ML  $\rightarrow$  we have to cluster our data.

## PCA

Refer "Application of Linear Algebra in Dimensionality Reduction"  $\rightarrow$  PREV. CHAPTER

Ex: #importing

```
from sklearn.datasets import load_breast_cancer
cancer_dataset = load_breast_cancer()
print(cancer_dataset.DESCR), cancer_dataset.keys()
df = pd.DataFrame(cancer_dataset['data'], columns=cancer_dataset['feature names'])
df.head()

from sklearn.preprocessing import StandardScaler # Standardization
scaler = StandardScaler()
scaler.fit(df)

scaled_data = scaler.transform(df)

# Applying PCA Algorithms
from sklearn.decomposition import PCA
pca = PCA(n_components=2)

data_pca = pca.fit_transform(scaled_data)
data_pca  $\rightarrow$  2 features
pca.explained_variance_
plt.scatter(data_pca[:, 0], data_pca[:, 1], c=cancer_dataset['target'])

plt.xlabel('First Principal Component')
plt.ylabel('Second Principal Component')

# helps in giving diff. color to points based on 'target' col. features
# ex: if Y/N 2 colors
```

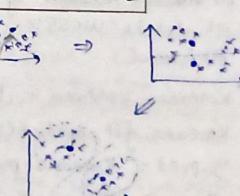
## K MEANS CLUSTERING UNSUPERVISED ML

Ex: 

Steps: 1) Initialize some 'k'  $\rightarrow$  centroids

2) Points that are nearest to the centroids  $\rightarrow$  group

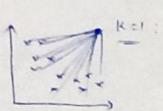
3) Move the centroids  $\rightarrow$  Avg. point will be the new centroid location.



- How do we select K value?

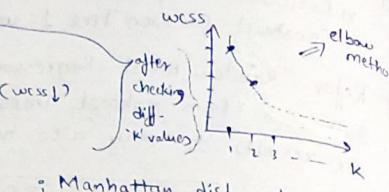
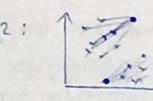
⇒ WCSS = within cluster sum of squares

⇒ Initialize K=1 to 20



$$WCSS = \sum_{i=1}^n (\text{distance b/w points to nearest centroid})^2$$

K=2:

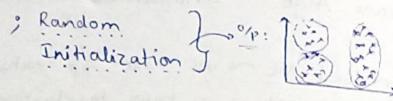
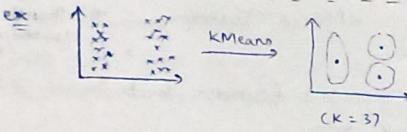


$$\Rightarrow \text{Euclidean dist.} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

(default)

$$\text{; Manhattan dist.} = |x_2 - x_1| + |y_2 - y_1|$$

- Random Initialization Trap:



To resolve this mismatch, we use "KMeans++ Initialization Technique".  
it initializes the centroid in such a way that it's completely at a farther distance between each other.

Ex: # Importing

```
from sklearn.datasets import make_blobs
```

```
x, y = make_blobs(n_samples=1000, centers=3, n_features=2)
```

```
x
```

```
y
```

```
plt.scatter(x[:,0], x[:,1], c=y)
```

```
①
```

```
②
```

```
from sklearn.cluster import KMeans
```

# Elbow method to select K value

```
wcss = []
```

```
for K in range(1, 11):
```

```
Kmeans = KMeans(n_clusters=k, init="k-means++")
```

```
Kmeans.fit(x_train_scaled)
```

```
wcss.append(Kmeans.inertia_)
```

```
wcss
```

```
# Plot elbow curve
```

```
plt.plot(range(1,11), wcss)
```

```
plt.xticks(range(1,11))
```

```
plt.xlabel("Number of Clusters")
```

```
plt.ylabel("WCSS")
```

```
plt.show()
```

```
Kmeans = KMeans(n_clusters=3, init="k-means++")
```

```
Kmeans.fit(x_train_scaled)
```

```
y_pred = Kmeans.predict(x_test_scaled)
```

```
plt.scatter(x_test[:,0], x_test[:,1], c=y_pred)
```

to decide cluster value  
from the graph

the point where there  
is not much change  
in 'y' is the ideal  
point.

To validate 'K' value, we have 2 more methods: kneelocator, silhouette scoring

From prev. graph -> convex ✓  
(since it's kind of convex, so curve convex)  
Knee Locator (range(1,11), wcss, curve = "convex", direction = "decreasing")  
Knee elbow → We can validate our elbow point from this  
(prev. curve is convex so we write, dir = "decreasing")

silhouette scoring

```
from sklearn.metrics import silhouette_score
```

```
silhouette_coefficients = []
```

```
for K in range(2, 11): # We start at 2 clusters for silhouette coefficients
```

```
Kmeans = KMeans(n_clusters=k, init="K-Means++")
```

```
Kmeans.fit(x_train_scaled)
```

```
score = silhouette_score(x_train_scaled, Kmeans.labels_)
```

```
silhouette_coefficients.append(score)
```

```
silhouette_coefficients
```

(highest peak no. of clusters)

```
# plotting silhouette score
```

```
plt.plot(range(2,11), silhouette_coefficients)
```

```
plt.xlabel("Number of Clusters")
```

```
plt.ylabel("Silhouette Coefficients")
```

```
plt.show()
```

sort Kmeans what label  
we have got i.e. to which  
cluster that specific datapoint  
belongs to.

Here we don't consider centroids

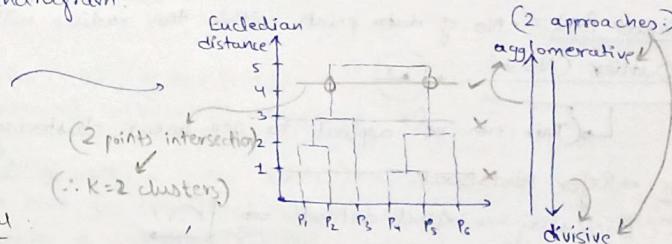
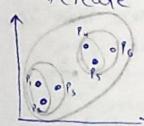
## HIERARCHICAL CLUSTERING

Steps: 1) Each point will be considered as a separate cluster.

2) find the nearest point & create a new cluster

3) Keep on doing the same process until we get a single cluster.

4) Create a dendrogram.



• K-Means vs Hierarchical

1) Dataset size → Huge = K-Means

Small = Hierarchical clustering

2) K-Means → Numerical Data

HC → variety of data

(can use cosine similarity)

3) K-Means = Elbow Method = no. of centroids.

HC = no. of clusters

from sklearn.preprocessing import StandardScaler  
scaler = StandardScaler()  
x\_scaled = scaler.fit\_transform(iris\_data)  
# Apply PCA = reducing the dimensions  
from sklearn.decomposition import PCA  
pca = PCA(n\_components=2)  
pca\_scaled = pca.fit\_transform(x\_scaled)  
plt.scatter(pca\_scaled[:,0], pca\_scaled[:,1], c=iris.target)

## # Agglomerative Clustering → constructing a dendrogram

import scipy.cluster.hierarchy as sc

plt.title("Dendrogram")

sc.dendrogram(sc.linkage(pca\_scaled, method='ward'))

plt.xlabel('Sample Index')

plt.ylabel('Euclidean Distance')

from sklearn.cluster import AgglomerativeClustering

cluster = AgglomerativeClustering(n\_clusters=2, metric='euclidean', linkage='ward')

cluster.fit(pca\_scaled)

cluster.labels\_

plt.scatter(pca\_scaled[:, 0], pca\_scaled[:, 1], c=cluster.labels\_)

## # Silhouette Scoring

from sklearn.metrics import silhouette\_score

silhouette\_coefficients = []

for k in range(2, 11):  
 (We start at 2 clusters for silhouette coeff.)

agglo = AgglomerativeClustering(n\_clusters=k, metric='euclidean', linkage='ward')

agglo.fit(x\_scaled)

score = silhouette\_score(x\_scaled, agglo.labels\_)

silhouette\_coefficients.append(score)

plt.plot(range(2, 11), silhouette\_coefficients)

plt.xlabel('Number of Clusters')

plt.ylabel('Silhouette Coefficient')

plt.show()

## DBSCAN CLUSTERING

→ Core Point: • No. of points within the ' $\epsilon$ ' should be greater  $\geq$  minpts.  
• Let minpts = 4,  $\epsilon$  = radius  
• (other datapoints)

→ Border Point: • No. of data points within this radius will be less than minpts.

→ Outlier (Noise):  
• (This can be applied to non-linear clustering as well.)

Ex: → Refer Notebook → (similar way written in SVM kernels)

Basically → initial dataset →

→ standardizing & applying DBSCAN. → [from sklearn.cluster import DBSCAN]

## SILHOUETTE CLUSTERING

Steps: 1) Create a specific cluster & then wrt a particular point calculate avg. distance.

$$a_i = \frac{1}{|C_{i-1}| - 1} \sum_{j \in C_{i-1}} d(i, j)$$

(No. of points belonging to cluster  $i$ )



2) We take another cluster & then from that same point, we try to calculate the avg. distance to the other points in the 2nd cluster.



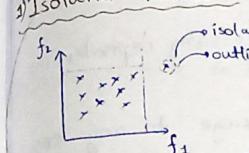
$$b_i = \min_j \frac{1}{|C_{j-1}| - 1} \sum_{k \in C_{j-1}} d(i, k)$$

$$\text{3) } s(i) = \begin{cases} 1 - a(i)/b(i), & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1, & \text{if } a(i) > b(i) \end{cases}$$

$-1 \leq s(i) \leq 1$   
↓  
(More near to +1, better clustering model we have created.)

## ANOMALY DETECTION MACHINE LEARNING ALGORITHMS

### 1) Isolation Forest:



Internally it creates many decision trees which are known as "isolated trees" wrt each datapoint.

$$\text{Anomaly Score: } S(x, m) = \frac{-E(h(x))}{c(m)}$$

- $m$  = no. of data points
- $x$  = data point

- $E(h(x)) = \text{average search depth for } x \text{ from the isolated tree.}$
- $c(m) = \text{average depth of } h(x)$

→  $E(h(x)) \ll c(m) \Rightarrow S(x, m) \approx 1 \Rightarrow \text{Anomaly Score} = \text{OUTLIERS}$   
→  $E(h(x)) \gg c(m) \Rightarrow S(x, m) \approx 0.5 \Rightarrow \text{Normal data point}$

Ex: df = pd.read\_csv("healthcare.csv")

plt.scatter(df.iloc[:, 0], df.iloc[:, 1])

from sklearn.ensemble import IsolationForest

clf = IsolationForest(contamination=0.2)

clf.fit(df)

predictions = clf.predict(df)

index = np.where(predictions < 0) → Outliers

index

x = df.values

plt.scatter(df.iloc[:, 0], df.iloc[:, 1])

plt.scatter(x[index, 0], x[index, 1], edgecolors='r')

	0	1
0	1.617	1.944
1	1.256	1.609
2	-2.349	4.392

### 2) DBSCAN Clustering:

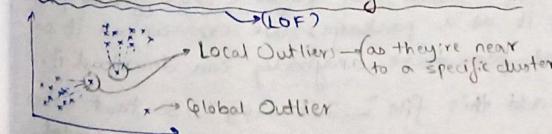
#### Advantages:

- No need for predefined cluster count (unlike K-means).
- Capable of clustering non-linear separable data.
- Robust to outliers: efficiently detects & excludes outliers.
- Few parameters needed (like  $\epsilon$  & min. pts)
- Fast region queries

#### Disadvantages:

- Border points can belong to multiple clusters, causing ambiguity.
- Sensitive to parameter selection.
- Dependency on distance measures (i.e. Euclidean/Manhattan).

### 3) Local Outlier Factor Anomaly Detection:



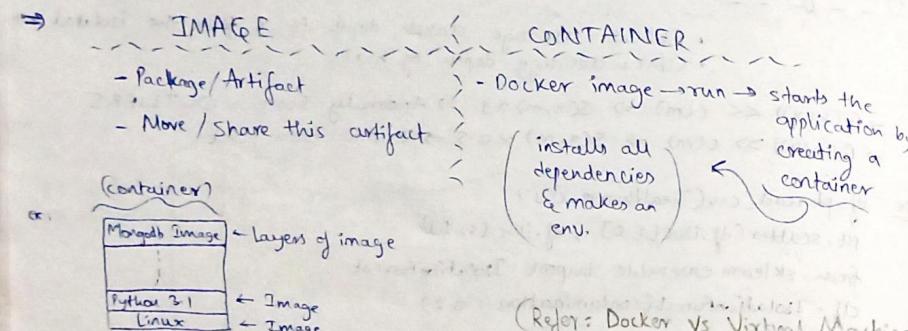
- LOF → uses KNN internally
- We will take first 'k' values from a point, then we take its K neighbours of those K points & we check the density of points

Locality is given by KNN, whose distance is used to estimate the local density. By comparing the local density of a sample to the local densities of its neighbors, one can identify samples that have a substantially lower density than their neighbors. These are considered outliers.

## DOCKER

### Containers:

- A way to package application with all the necessary dependencies & configuration.
- Portable artifact; easily share & move this package to any env.
- Makes development & deployment more easy & efficient.



Docker - Compose: a tool that helps us define & share multi-container applications.

(docker-compose.yml)

Volumes: preferred mechanism for persisting data generated by & used by docker containers.

## PROJECTS

### Project 1:

- Created a repository on Github.
  - In VS code terminal: → conda create -p venv python==3.8 -y → Creates venv folder
  - (conda activate venv)
- Create .gitignore file.
- Create setup.py, requirements.txt → (pip install -r requirements.txt)
  - setup.py = responsible for creating my ML application as a package. Now we can also install this package in different projects. After making it as a package, we can upload it on Python PyPi and from there anybody can download it.
- Create SRC folder & add this file "\_\_init\_\_.py", so that when find\_packages() from setup.py searches all the packages needed;

it will consider the SRC folder as a package.

Write setup.py

In requirements.txt, adding '-e', so that it automatically triggers setup.py.

Create components folder & create \_\_init\_\_.py file.

↳ (modules that we are going to specifically use in the project)

Creating these files inside components folder:

- data-ingestion.py
- data-transformation.py
- model-trainer.py

Create pipeline folder. Create these files in this folder:

- \_\_init\_\_.py
- train-pipeline.py
- predict-pipeline.py

Create logger.py, exception.py, utility in src folder. Write own custom exception.py. Write logger.py to log the exceptions which has occurred.

Create notebook folder & refer the dataset & both the .ipynb files.

Write → data-ingestion.py → data-transformation.py → model-trainer.py

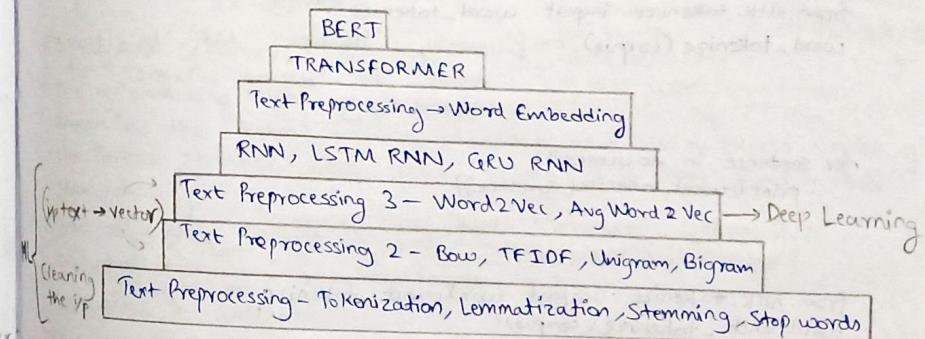
Write app.py & predict-pipeline.py

Deploying on AWS, Azure.

\* Refer: MLflow, BentoML, DugHub !

## NLP FOR MACHINE LEARNING

### Overview:



### Tokenization:

CORPUS = Paragraph      DOCUMENTS = Sentences

VOCABULARY = Unique Words      WORDS = words

Tokenization: Process wherein we convert paragraphs into sentences & sentences into words.

(aka Tokens)