

UNIT-3

Page No. _____
Date _____

~~df = pd.Series~~

```
df = {'a':100, 'b':200, 'c':300, 'd':400, 'e':500}  
print('Original', df)  
new_se = pd.Series(df)  
print('New Series', new_se)
```

Jupyter, Python, R

Jupyter

Jupyter notebook is an open source, web based interactive environment which allows to create and share documents that contain live ~~part~~ code, mathematical eqⁿ, graphs, maps, plots, visualizations, and narrative text. It integrates with many programming languages like Python, PHP, R, C# etc.

Advantages

1. All in one place - Since Jupyter is an open source web based interactive environment that combines code, text, Images, videos, mathematical eqⁿ, plots, maps, GUI and widgets to a single document.
2. Easy to convert - Jupyter notebook allow users to convert the notebook into other formats such as HTML & PDF. It also has online tools and which allows you to render a publically available notebook in the browser directly.
3. Easy to Share - JN can be shared in JSON format which makes them easily sharable

1 Language Independent :- It is platform independent because it is represented as JSON format which is a language independent text based file format, it can also be converted to any file format such as markdown down, HTML, PDF etc.

2 Interactive Code - It uses Jupyter widgets package which provides many common user interface for exploring code and data automatically.

Disadvantages

1. It is very hard to test long asynchronous tasks.
2. ~~less security~~
3. It runs cell out of order. code correction.
4. No IDE Integration, No Linting and no code style correction.

Components of Jupyter Notebook

1. Notebook Web Application - It allows users to edit code in the browser with automatic syntax highlighting & indentation.
- * Run Code on the browser
- * See results of computation with media representation such as HTML, PNG, PDF, LaTeX etc.
- * Create and use JS widgets
2. Kernels - Kernels are the separate processes started by the notebook web application to run a user's code in the given language and return output to the notebook web app.
In JN Kernel is available in following lang. Python, Julia, Ruby, R, Scala, Go, Node.js.

3 Notebook Documents - contains a rep. of all contents which is visible in the notebook web application, including input and output computation, text, mathematical eqn, graph, images.

Types of cells in JN

1. Code cell

2. Markdown cell → for documentation

① Bold & Italic (font styling)

bold **italic**

② Headers

H₁ → H₆

Heading 1 ## Heading 2

<h1> Header 1 </h1>

<h2> Header 2 </h2>

<h3> Header 3 </h3>

<h4> Header 4 </h4>

<h5> Header 5 </h5>

<h6> Header 6 </h6>

③ Ordered list

↳ 1 A

1 A 1

2 A 2

2 B

1 B 1

2 B 2

④ Bullet list

→ *

- A

* A.

- B

* B.

A.

B.

Hyper links

[Google] (www.google.com)

Table content

Markdown cell allows you to create a table using pipe symbol (|) and dash symbol (-)

| → column | Name |

— → rows | --- |

Ex:- | Round | Name | Age |

= | --- | --- | --- |

| 01 | Jay | 18 |

| 02 | Amit | 19 |

Round Name Age

01 Jay 18

02 Amit 19

Images

To insert a image in a markdown cell you first need to insert the image in the same directory and for this go to Jupyter dashboard and select upload, and then specify the path of an image and click on open. after that go to your current notebook and type the following code ~~img src = "image.png"~~

3. Raw NB convert cell

It provides a place where you can write output directly. These cells are not evaluated by the notebook Kernel.

4. Heading Cell

The JN does not support heading cell and Jupyter shows the pop-up.

Ques Diff b/w Python, Jupyter, Python

Ans list out any 5 companies that integrate with ipython and for Jupyter

Ques what tools integrate with ipython & with Jupyter

Ans 1 ipython

- ipython is an interactive command-line terminal for Python.

- It is an enhanced version of Python shell.

Ans 2 python

Python is a general-purpose programming language i.e. widely used for web dev. AI, DS.

Ans 3 Jupyter

Jupyter is a web-based interactive computing environment that supports several prog. lang. and share doc containing code, e.g., visualization text etc.

Any 2
 Jupyter - Trivago, Roff, Frappe, monolith 4.2, Explor
 Jupyter - MS Azure, Google Colaboratory, Uber, Amazon
 Amazon Sage Maker, KBR, Nscorp

Any 3
 Tools - Jupyter → Jupyter contains notebooks,
 Magic functions, widgets, Nbviewer, nbqa.

Ipython - matplotlib, numpy, pandas, scipy, seaborn

Understanding the Tools

↳ current session details

Magic commands

%quickref to tell how to obtain additional
 help info. about -

- object? ↳ 1. Using ipython to perform calculations
- ? ↳ 2. Obtaining info. about magic fun
- 3. Learning about python prog. language
- 4. Discovering facts about the packages,
 objects, methods that you use in python
 to interact with data.

Ipython provides a standard feature called as
 cls (clear screen) which is not present in py console

Changing window Appearance

1. Option:

a. Font:

3. Layout:

4. Colors: %colors

Getting Python Help

modules → Help mode → full details

→ Interactive Help Mode → Particular detail of obj.

We use help mode when you want to explore the lang.

When we specifically need help with an object

There are certain commands to obtain a listing of some other topics

Modules - current session list

Keywords - list of all keywords

Topics -

Using Python object help

If you have a list called mylist then just by typing mylist? we can discover obj type, content, length.

Ques :- what do you understand by Raw NB convert

```
%pdoc
%pdf
%source
%file
%pinfo
%pinfo2
```

Performing multimedia & Graphic Integration

→ Loading from online sites

% load http:// _____ .py

→ Obtaining online graphics and multimedia

* Ipython.display

from Ipython.display import Image

Embed = Image(' ', '')

Embed

or

Softlinked = Image(url=' _____ ')

↳ If the image is ~~referred~~ changed
from the website then it will change
automatically.

To change image format

from Ipython.display import set matplotlib formats
set_matplotlib_formats(['pdf', 'svg'])

To convert a ~~csv~~

To create a csv file from Dataframe

Load & read csv

Dataframe to csv

Import pandas as pd

```
df = pd.DataFrame([[1, 2, 3, 4], [5, 6, 7, 8], [9, 10, 11, 12]], columns=['A1', 'B1', 'C1', 'D1'])
```

df.head()

df.to_csv('export.csv')

or

df.to_csv('export.csv', index=False)

A1 B1 C1 D1

0 1 2 3 4

1 5 6 2 8

2 9 10 4 12

Pandas → concat

```
dict1 = {'Name': ['A1', 'B1', 'C1'],
         'Age': [12, 18, 20],
         'Address': ['xyz1', 'abc2', 'def3'],
         'Qualification': ['10th', '11th', '12th']}
```

}

dict2 = {'Name':

}

df = pd.DataFrame(dict1, index=[0, 1, 2, 3])

index = 0, 1, 2, 3

index = 0, 1, 2, 3

df2 = pd.DataFrame(dict2, index=[4, 5, 6, 7])

frames = [df, df2]

res1 = pd.concat(frames)

print(res1)

use of Matplotlib

```
from matplotlib import pyplot as plt
plt.plot([1, 2, 3], [3, 5, 7])
plt.show()
```

```
plt.title('Today Lab')
plt.ylabel('y-axis')
plt.xlabel('x-axis')
```

Ques.

- Draw a graph containing comparison of lines in a single plot and also perform the following.
- ① change the style of the plot.
ggplot()
 - ② Add the labels on the graph using legend()
legend function

```
from matplotlib import pyplot as plt
import numpy as np
from matplotlib import style
data = np.random.randn(50)
plt.style.use('ggplot')
plt.plot(data, linestyle = ':', linewidth = 2)
```



Page No. _____
Date _____

Magic Commands :- are special commands that can help you with running & analysing data in your notebook.

They add a special functionality i.e. not straight forward to achieve with python code or Jupyter notebook interface.

1. Line Magic - It is rep. using % prefix and operates on a particular line of input.
It is used in the form of exp. and their return value can be assigned to a variable.

2. Cell Magic - It is rep. using %% prefix and works on a complete cell or multiple lines of input.
They receive the whole block as a string.

Built-In Magic Commands

1. %autocall :- make func callable without having to type parentheses
2. %automagic
3. %automagic
4. %cd
5. %dhist
6. %edit
7. %env
8. %gui
9. %lpmagic
10. %notebook
11. %who
12. %pwd
13. %recall
14. %matplotlib -- list
15. %time
16. %quickref
17. %load

Magic Functions

1 %autocall :- Make func callable without having to type parentheses.

Ex: `def add(x, y):
 return x+y`

%autocall

`add 3,5`

→ Output

→ `add(3,5)`

→ `add(3,5)`

8

2 %automagic :- Make magic func callable without having to type the initial %.

Ex:-

%automagic

`cd, hist, ls etc.`

3

%cd :- change the current working directory
%cd ..

4

%hist :- Print your history of visited directories

5

%env :- Give a list of environment variables

6

%run :- Run the named file inside Jupyter as a program
%edit :- Open editor to execute code with history
of current session.

Ex:-

`run hello.py`

`edit hello.py`

7

%lsmagic :- List currently available magic functions

8

%matplotlib --list :- Available matplotlib backends
['ttk', 'gtk', 'svg', 'pdf'] etc

9

%who :- Print all interactive variables, with some minimal formatting.

10

%pwd :- Return the current working directory path

To draw figures → %matplotlib inline
%matplotlib

Page No.	
Date	

1. load

11. load :- Load code into current frontend file.py
12. %quickref :- Show a quick reference sheet
13. %recall :- Repeat a command, or get a command to input line for editing Eg:- %recall 45^{line no.}
14. %notebook :- This fun^ct converts current Jupyter history into an Ipython notebook file with ipynb exten^{sion}.
Eg:- %notebook filename

15. %time :- Time execution of a python statement or expression.

[Data cleaning]

1 Empty cells → dropna()

df = pd.read_csv('data.csv')

① df.dropna(inplace=True)
To remove from original data

② df.dropna()

③ print(df.to_string())

④ fillna

0 or mean, median, mode.

⑤ df = pd.read_csv('data.csv')

df.fillna(130, inplace=True)

⑥ x = df['column'].median()

~~df['column'].fillna(x, inplace=True)~~

Data of Wrong Format

To date (Wrong format of date)

```
df['date'] = pd.to_datetime(df['date'])
print(df.to_string())
```

```
df = pd.read_csv('data.csv')
for x in df.index:
```

```
if df.loc[x, 'Duration'] > 120
    df.loc[x, 'Duration'] = 120
```

* (df.duplicate())

X X X X X X

Working with real Data

Files 1. colors.txt , 3. values.xls
2. titanic.csv , 4. XMLData.xml

→ We have to access data in a no. of forms and locations

Eg :- Memory streams rep. a form of data storage that your computer support natively, flat file exist on your hard drive, relational databases usually appear on networks. (although smaller relational data bases like that of MS access on a H.D), and web based data from the internet

Uploading, Streaming & Sampling Data

DS call the columns in a database features/variables.

The rows are called as cases. Each row rep. a collection of variables that you can analyse.

Uploading

1. Small amount

color.txt

color	val
Red	1
Blue	2

with open('color.txt', 'r') as open_file:

```
print("Color.txt Content: " + open_file.read())
```

In this the entire data set is loaded from the memory into free memory.

Streaming

Some data sets are usually large such that you wouldn't be able to fit them entirely in memory at one time. There are also cases when some data sets load slowly bcoz they reside on a remote site.

Streaming answers both needs by making it possible to work with the data a little at a time.

Note:- You download individual pieces, making it possible to work with just part of the data and to work with it as you receive it, rather than the entire data set to download.

with `open('color.txt', 'r')` as open-file:
 for observation in open-file:
`print('color.txt content: \n' + observation)`

* Python streams each records from the source.
 This means that you must perform a read
 for each record. You want

#

Sampling

Data streaming obtains all the records from a data source but sometimes you don't need all the records, and in such case you can save time & resources by sampling the data. This means retrieving records a set no. of records apart such as every 5th record or by making random samples.

$n = 2$

with `open('color.txt', 'r')` as open-file:

for j, observation in enumerate(open-file):

if $j \% n == 0$

`print('Reading line: ' + str(j) + ' content: \n' + observation)`

* enumerate fun is use to retrieve a row no. When $j \% n == 0$ the row will be kept and the application outputs that row. Otherwise it will discard the value.

Transparent

read_table → reads tabular data
and create a dataframe

Parsing → Breaking down of elements

Parsing → Discarding codes into blocks
Interpreting the code on data

Date	No.

Values b/w → 0 and 1

Random Sampling

```
from random import random
sample_size = 0.25
with open('color1.txt', 'rb') as open_file:
    for j, observation in enumerate(open_file):
        if random() <= sample_size:
            print('Reading line:' + str(j) +
                  ' content:' + observation)
```

Accessing Data in Structured flat file form

- A flat file presents the easiest kind of file to work with the data appears as a simple list of entries that you can read one at a time in memory.
- Classes and methods in the Pandas library parses the flat file data to make it easier to manipulate.
- A text file generates treats or data as strings so that one needs to convert numeric data in other forms.
- A CSV files provide more formatting and more info. but it requires more effort to read this file.

Parag 2.

Reading from a text file

import pandas as pd

color_table = pd.read_csv('color1.txt')

print(color_table)

Within Pandas there is a set of parsers, code used to read individual bits of data and determine

The purpose of each bit acc. to the format
of the entire file.

Notes
Pandas interprets the 1st row as field names.

- d. Reading CSV delimiter file
- i. Header, fields, records, strings, integers, real no.

titanic = pd.io.parsers.read_table('titanic.csv')

X = titanic[['age']]

print(X)

titanic = pd.io.parsers.read_table('titanic.csv', sep=',')

↳

- # Reading an Excel and other MS applications - highly formatted content. One can specify every aspect of info. these files contains.

Xls = pd.ExcelFile('value.xls')

figvalues = Xls.parse('Sheet 1', index_col='name', na_values=['na'])

print('fig values')

Sending Data in an Unstructured File format

- Unstructured data files consist of series of bits
- The file does not separate the bits from each other in any way.
- Unstructured file format rely on the file user to know how to interpret the data.

Eg:- Each pixel of a picture file could consist of 3

fields of 32 bits each showing this depends upon the data scientist. However the beginning of the file may provide clues about interpreting the file.

Install → pip install ~~skimage~~ scikit-image

The imshow() function of matplotlib performs the rendering and uses a grey scale color map.

The show() func. is used to just display the image for you.

- # The image is actually an array of pixels that can be manipulated in different ways for example cropping of image (the purpose of cropping the image is to make it of a specific size both images must be of the same size for analysing purpose). Cropping is one way to ensure that the images are the correct size for analysis.
- Another way to change the image size is to resize it.
- Once you have the images of right size you need to flatten them. A dataset now is always a single dimension, not 2-dim.
- The type of this array is numpy.ndarray you can add this array to a data set and then use that data set for analysis purpose.

Managing Data from Relational Database

The goal of db manager is to make data easy to manipulate the focus of most data storage is to make data easy to retrieve.

SQL shapes data in a particular way so as to ease & eliminate the shaping procedure (creating a connection to a db is a complex procedure)

Once you have access to an engine you can use that engine to perform task specific to a DBMS. The output of a read method is always a dataframe object that contains the requested data.

1. `read_sql_table()` - Reads data from SQL table to ~~read_sql()~~ dataframe object.
2. `read_sql_query()` - Read data from a database using a SQL query to a dataframe object
3. `read_sql()` - Read data either from a table/query
4. `Dataframe.to_sql()` - write the content of a ~~Dataframe~~ object to the specified table in the database

SQLAlchemy

\hookrightarrow SQLite, MySQL, PostgreSQL, SQL Server, ODBC

PostgreSQL

Interacting with data from NoSQL db
 MongoDB → PyMongo, MongoClient
 This class is used to create the required engine.

The MongoDB engine relies on the find() function to locate the data.

Accessing Data From the Web

Other service is a kind of web application that provides a means to ask questions and receive answers. Web services usually host a no. of input types.

There is a service called as micro-service (query system) that has a specific focus and provide only one specific query input and output.

Working with Web services and microservices means working with ~~XML~~ XML.

<MyDataset>

<Record>

<Number> | </Number>

<String> Welcome </String>

<Boolean> True </Boolean>

</Record>

{

</MyDataset>

Visualisation

Page No. _____

Date _____

- # Performing data analysis is basically about the communication because unless you communicate your idea, the act of obtaining, shaping and analysing data has a little value beyond individual personal needs.
- Matplotlib in cmp to matlab is relatively easy bcoz they use same sort of state machine to perform fast and they have similar method of defining graphical elements.

Starting with Graph

A Graph or chart is simply a visual rep. of numeric data.

1. Defining the plot.
2. Defining multiple lines to plot
3. Saving your work

Setting the Axis Ticks, Grids

- The use of Axes, Ticks & grids make it possible to elaborate graphically the relative size of data elements so that the viewer gains an appreciation of comparative measure

1. Setting the axis
2. Formatting the axis
3. Adding grid.

Matplotlib

Matplotlib is a low level graph plotting library that serves as a visualization utility.

1. Matplotlib Version

2. Matplotlib Installation

3. Use of Pyplot sub modules

4. Plotting

a. Use of plot function

- Plotting without line
- Multiple points
- Default 2 points

5. Matplotlib Markers

- Use of marker size ('ms')
- Marker color ('mec')
- Marker face color ('mfc')
- Format strings ('fmt') (~~for~~
(Marker/Line/Color))

6. Matplotlib Line

- Types of line reference (Default dotted, dashed, dashdot)
- Line style ('ls') . Line color ('c')
- Line width ('lw') . Multiple lines

7. Matplotlib Labels and titles

- xlabel, ylabel, title, font properties for title and label
(using fontdict parameter)

- Positioning the title ('loc')

8. Add grid lines

a. using grid function

- x-axis, y-axis grid, 3. color of grid, 4. line properties (color, line style, line width)

Subplots()

To draw multiple plots in 1 figure.

Add title to each plot

supertitle

Practice Questions - 1

Q1 ① dict = {

'Name': ['Raman', 'Sheetal', 'Vikas', 'Shiv'],
 'Age': [18, 17, 15, 20],
 'Sex': ['Male', 'Female', 'Male', 'Female'])

import pandas as pd
 import numpy as np

② df = pd.DataFrame(dict)

df.columns

Output ['Name', 'Age', 'Sex']

③ Sort the df based on numerical column

df['Age'].sort() ascending = False

df.sort_values('Age')

④ Mean

men = df['Age'].mean()

⑤ df['Name'].runa(men)

⑥ `nl = df.columns[df.isnull().sum() > 4]`
`df.dropna(nl)`

	col-A	col-B	
0	A	1	→ date.csv
1	B	2	
2	C	3	
3	D	4	

`df.insert(2, 'col-C', Data)`

`idx = 1`

`val = [2, 4, 6, 8]`

`df.insert(loc=idx, column='col-C', value=val)`

Output	col-A	col-C	col-B
0	A	2	2
1	B	4	2
2	C	6	3
3	D	8	4

⑦ `df.select_dtypes(include=['Bool']).columns`

⑧ `df.count()` → No. of non-Nan values in each column

`df.sum(df.count())` → sum of total count values which are not NAN

⑨ `df=df-split = np.array-split(df, 2)`
 for df in df-split:
`print(df)`

⑩ Reversing Row wise

~~df[1].iloc[:,-1]~~
~~df[1].iloc[:,2]~~
~~df[1].iloc[:,1]~~
~~df[1].iloc[:,0]~~

Rewriting column wise
`df[df.columns[::-1]]`
~~df[1].iloc[:,0]~~
~~df[1].iloc[:,1]~~
~~df[1].iloc[:,2]~~

Q2

Purpose

Python
General Purpose programming language

Syntax

uses indentation and is more readable.

Data types

Supports dynamic typing and a wide range of data types.

Libraries

large ecosystem of libraries and packages are available

Performance

Interpreted language, which can be slower for certain tasks

Community support

large community support

Applications

Web development, data analysis, ML etc.

Q2

Purpose

3D visualization

Integration

Matplotlib
Python library for data visualization

Supports 3D plotting and visualization

Can be used with other Python libraries and tools

MATLAB

Numerical Computing and technical computing environment

uses traditional programming syntax

Provides built-in support for numerical and matrix operations.

Rich set of toolboxes for various scientific applications.

Compiled language generally faster for numerical computation.

smaller community support

Mathematical, engineering and scientific research

MATLAB

- Integrated plotting and visualization in MATLAB

- Offers built-in 3D plotting capabilities.

- Designed to work seamlessly within MATLAB.

uses Python syntax

open-source and free to use.

Data analysis, scientific plotting, visualization.

`X.flatten()`

Returns a copy of the array as a flattened 1D array.

Allocates \rightarrow returns a new array.

May be slower due to memory allocation.

Does not modify the original array.

Returns a 1D array regardless of the original shape.

uses MATLAB syntax

Commercial software with license fees.

Mathematics, engineering, and scientific research.

`X.ravel()`

- Returns a view (if possible) or a copy of the array.

- Returns a view whenever possible, avoid memory allocation.

- Can be faster since it avoids memory allocation.

- Does not modify the original array.

- Preserves the shape of the original array.

Note # The "flatten()" fun always returns a copy of the array as a flattened 1D array, regardless of the original shape.

- It guarantees a new memory allocation and returns a contiguous array.

- The 'ravel()' returns a view of the original array whenever possible, avoiding memory allocation.

- It preserves the order of elements as stored in memory and maintains the shape of the original array.

84

fun.

usage.

control.

plots
Affected

Returns.

multiple
subplotsex:-
=**axis**

Controls the units and scaling of the plot's data.
Applied to an existing plot or subplot.

Controls the axis units, labels and tick marks.

Affects the entire plot or subplot.

Returns the current axis units and scaling.
Affects all subplots if applied to the entire figure.

`plt.axis([xmin, xman, ymin, ymax])`

axes

- Represents the plotting area on a single subplot.
- Used to create and manipulate plots & subplots.
- Provides methods to add and customize plot elements.
- Operates on a specific subplot or plotting area.
- Returns a reference to the created 'axes' object.
- Associated with a specific subplot within the figure.

`fig, ax = plt.subplots()`

Note
=

Axes is the region where the plots and graphs are drawn.

Each axes has its own title, x-axis and y-

Axis in matplotlib refers to the scale of the graph. A two-D graph has an x-axis and y-axis, whereas a 3-D graph has x, y and z-axis.

- # 1. Defining the line appearance
2. Working with line style
3. Using colors
4. Adding markers

Using Label, Annotation and Legends

Label - Its purpose is to make it easy for the viewer to know the name or kind of data that is to be illustrated.

Annotation - It augments the info. the viewer can immediately see about the data with notes, sources or any other useful info.

In contrast to label the purpose of annotation is to help extends the viewer's knowledge of the data rather than simply identify it.

Legends

Presents a listing of the data groups within the graph and often provides clues to make identification of the data group easier.

Practice Question - d

Q1
= Read total profit of all months and show in a chart
using plot.
plt.plot(total-profit)

Q2
Legend styling

Q1
= Import pandas from pd
Import numpy as np
Import matplotlib.pyplot as plt

df = pd.read('---')

Q1
= plt.plot(df['total-profit']), plt.title('Total Profit'),
plt.show()

Q2
= plt.plot(df['total-profit'], c='red', ls='--', marker=1)
plt.title('Total Profit')

plt.plot(df['total-profit'], 'o:r', c='red', marker=1, mfc='g', lw=3)
plt.xlabel('Month Number')
plt.ylabel('Sold Units Number')
plt.legend(loc=4)

Q3
=

Q1
 $\text{toothpasteSales} = \text{df}[\text{'toothpaste'}]$
 $\text{mouthNo} = \text{df}[\text{'mouth_number'}]$
 $\text{plt.scatter}(\text{toothpasteSales}, \text{mouthNo})$
 $\text{plt.xlabel}(\text{"Toothpaste Sales"})$
 $\text{plt.ylabel}(\text{"mouth-No."})$
 $\text{plt.title}(\text{"Toothpaste Sales Data"})$
 $\text{plt.show}()$

Q5
 $\text{plt.bar}(\text{df}[\text{'mouth-no'}], \text{df}[\text{'face-cream'}])$
 $\text{plt.show}()$ $\text{plt.xlabel}(\text{'mouth'})$
 $\text{plt.ylabel}(\text{'Face cream'})$

Q6
 ~~plt.plot~~ plt.bar
 $\text{plt.pie}(\text{df}[\text{'Bathing-soap'}], \text{autopct} = \text{'%1.1f%'})$
 $\text{plt.savefig}(\text{'Bathing-soap.png'})$
 $\text{plt.show}()$

Q7
 $\text{plt.hist}(\text{df}[\text{'total-profit'}], \text{bins} = 10, \text{edgecolor} = \text{'red'})$
 $\text{plt.show}()$

Q8
 $\text{TotalSales} = \text{df}[\text{'total-units'}].\text{sum}()$
 $\text{plt.pie}(\text{TotalSales}, \text{autopct} = \text{'%1.1f%'})$
 $\text{plt.show}()$

Q9
 $\text{plt.subplot}(2, 1, 1)$
 $\text{plt.plot}(\text{df}[\text{'Bathing soap'}],$
 $\text{df}[\text{'month'}])$

- Q1 What is Faker Library? What it is use for. General
a dataset using Faker Library to create a dataset
of 100 fake names and email addresses.
- Q2 How to plot multiple bar charts using matplotlib
in python.
- Q3 Explain box plot in statistics

Ans 1

Faker Library generates fake (but reasonable) data for you which can be used for various testing and building applications. It is commonly used in software development testing and data analysis to generate realistic-looking but fictional data.

Usage

It is use to generate fake data in a wide range of categories, including names, addresses, phone numbers, email, addresses, job titles, credit card numbers, data and more. You can customize the data generation to match specific requirements or generate random data by default.

```
From faker import Faker
fake = Faker()
dataset = []
for _ in range(10):
    name = fake.name()
    email = fake.email()
    dataset.append([name, email])
for name, email in dataset:
    print('Name:', name)
    print('Email:', email)
```

A box plot, also known as a box-and-whisker plot, is a statistical visualization tool that provides a concise summary of data set's distribution. It displays the median, quartiles, and potential outliers of the data in a graphical form. The plot consists of a rectangular box and two lines, which resembles a box with whiskers.

import numpy as np

$x_1 = [2.5, 3.5, 4.5, 5.5, 6.5, 7.5]$

$y_1 = [3, 6, 9, 12, 15, 18]$

$x_2 = [5, 10, 15, 20, 25, 30]$

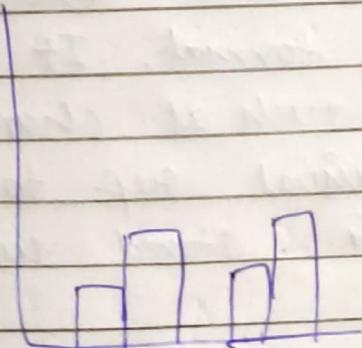
$y_2 = [2.5, 9.5, 14, 12, 8.5, 18]$

$r = np.arange(1)$

$width = 0.9$

$plt.bar(9, r, width=width)$

$plt.bar(x_1 + width, y_1, width=width)$



Ans 2

- Q1. write a short note on diff. customization option available with any plot.
- Q2. Explain Data Visualization
- Q3. With an Example of data comparison where we can use the scatter plot.
- Q4. What is open data. Name any website from which we can download open data.

Ans 1

- Plotting Libraries such as Matplotlib and Seaborn offers various customization options to enhance the visual rep. of data.
1. Color and Style
 2. Axis Labels and Titles
 3. Legends and Annotations
 4. Plot Size and Layout
 5. Axis Ticks and Tick Labels
 6. Grid Lines
 7. Plot Annotations

Ans 2

2. DV in DS refers to the process of exp. and displaying the structure and content of a database in a visual format. It involves using visual elements such as chart, graphs, diagrams and other visual rep. to gain insights into the data stored within a database.

Matplotlib :- A plotting library that provides a MATLAB-like interface for creating a wide range of static, animated and interactive visualizations.

- + Features
- + Easy to use. → Versatile API → customization

Limitations

- + Limited 3D plotting
- Lack of interactivity

Q3 Import matplotlib.pyplot as plt

heights = [53, 68, 20, 61, 72, 16, 54, 69, 62, 21]

weights = [150, 180, 200, 135, 200, 160, 145, 175, 185, 190]

plt.scatter(heights, weights)

plt.xlabel('Height')

plt.ylabel('Weight')

plt.title('Height vs Weight')

plt.show()

Q4 Open data refers to data that is freely

available for anyone to access, use and share provided by govt, org with intention of promoting transparency.

Eg - data.gov, Kaggle, open data portal by European Union.