# INDEX

# INTRODUCTION :-

→ Introduction to Data Analysis, Preparing data for analysis is often the large portion of the job time-wine for those working with data.

→ It should be noted that data wrangling isn't something we do merely once in our analysis

→ It is highly likely that we will do some data wrangling and move on to another analysis task. Such as data Visualization.

→ We need our data to be in-fer the Visualization that best convey what we are looking to show, and the data points we should collect for our analysis.

→ This comes with experience, So we must practice the skills that are covered in this topic.

## Data Wrangling:

→ Data Wrangling is the process of removing errors and combining complex-

data sets to make them more accessible & easier to analyze.

→ Due to the rapid expansion of the amount of data & data sources available today, Storing & Organizing large quantities of data for analysis is becoming increasingly necessary.

→ A data wrangling process also known as a data munging process, Consists of Reorganizing, transforming and mapping data from one "raw" form into another in order to make it more usable & valuable for a variety of downstream uses including analytics.

## IMPORTANCE of DATA WRANGLING.

→ Data professionals spend almost 80% of their time Wrangling the data, leaving a mere 20% for exploration & modeling.

→ A simple analogy will help you understand. The foundation of a Skyscraper is expensive & time consuming before the above ground structure starts.



Fig:1: Tasks Of Data Wrangling.

(Circle labels: ENRICHMENT, STRUCTURING, CLEANING, PUBLISHING, DISCOVERING, VALIDATING — Task Of Data Wrangling.)

→ Similarly, for data handling, Once the Code and infrastructure foundation are gathered, it will deliver immediate results [Sometimes almost-instantly] for as long as process is relevant.

→ However stopping necessary data wrangling steps will lead to significant downfall, missed opportunities and erroneus models that damage the reputation of analysis unit on organization.

## Primary importance of data wrangling:

- Making raw data usable
- Getting all the data from various sources
- Piecing together raw data according to the required format.
- Automated data integration tools are used as data wrangling techniques.
- Cleaning the data from the noise con flawed, missing elements

- Data Wrangling acts as a preparation stage of the data moving Process.
- Helping business user make Concrete, timely decisions.

## TASK OF DATA WRANGLING:

- To prepare your data for analysis as part of data munging. there are 6 basic steps one need to follow,

They are

1. Data Discovery: This is an all encompassing term that describes understanding what your data is all about. You get familiar with your data.

2. Data Structuring: When you collect-raw data, it initially is in all shapes & sizes, and has no definite structure. Such data need to be restructured to suit that your enterprise plans to deploy.

## 3. Data Cleaning:

- Raw data comes with some errors that need to be fixed before data is passed on to the next stage.
- Cleaning involves the tackling of outliers, making Corrections con deleting bad data Completely.

## 4. Data Enriching:

By this stage, you have kind of become familiar with the data in hand.

Enriching the data is an optional step that you only need to take if your current data doesnot meet your requirements.

## 5. Data Validating:

Validating the data is an activity that services any issues in the quality of your data so they can be

- addressed with the appropriate transformations.

## Rules of Data Validation:

- Quality
- Consistency
- Accuracy
- Security
- Authenticity

## 6. Data Publishing:

→ By this time, all the steps are completed and the data is ready for analytics. All that's left is publish the newly wrangled data in a place where it can be easily accessed and used by you & other stakeholders.

→ The final output of your efforts will be high quality of data that you can use to gain insights.



Fig-2: Importance of Data Wrangling.

## Data Wrangling Tools:

1. Parsehub.
2. Scrapy.
3. Talend
4. Alteryx APA Platform.
5. Altair Monarch.
6. Microsoft Power Query.
7. Tableau Desctop.

## Introduction to Python:

- Python is a programming language that let you work quickly and integrate system more efficiently.

- There are two major Python version: Python 2 and Python 3. Both are quite different.

- It was designed with an emphasis on code readability.

- This programming language were quickly & integrate system more efficiently.

③

**Wondows:** There are many interpreters avaible freely to run Python Scripts like IDLE that Comes with Python Software downloaded from http://Python.org/

**Lonus:** Python Comes Preinstalled with popular Linux Such as Ubuntu and Fedora.

**Mac OS:** Python 2.7 Comes bundled with Mac OS & Manually Install Python 3 from http://Python.Org/

**Note:**
Python doesnot require a Semicolon at the end of each Statement. Comment begin with a '#' and extend to the end of the line.

**Python Source Code:**
Source file use the ".py" extensions and are Called "Modules" with module. hello.py Casiest way to run the shell Command.

---

Data Meant to be Read By Machines.

Data Can be stored on Many ways of format & typon. Some formats store data snaway Carily handled by Machines. While otzen store data on a way Carily readable by a human.

Mocrosoft Word document are an example of the latter. While CSV, JSON & XML are examples of the former.

- Comma-Separate Values [CSV]
- JavaScript Object Notation [JSON]
- Extensible Markup Lanuge. (XML)

**Understanding the Database Schema**

→ In this topic a Very Small database for a Social Media application. The data Consist of four table:

---

1. Users  2. Posts  3. Command
4. Likes.

A high level diagram of the database Schema is Shown below



→ Both Users & Posts will have One to many relationship. III One Users Can post many Comments.
→ One post Can post also have multiple Comments.

④

# UNIT II

## INTRODUCTION TO PYTHON SQL LIBRARIES

→ All Software application interact with data, most common through a database Management System [DBMS]

→ In this unit, you will explore the different Python SQL Libraries that you can use & Develop a application to Interact with SQLite, MySQL & PostgreSQL database.

- Connect to different DBMS with Python SQL Libraries.
- Interact with SQLite, MySQL & PostgreSQL
- Perform common database queries using Python application
- Develop application across different database using a Python Script.

### SQLite:

SQLite database are Serverless and Self-Contained, Since they need to write data files. This means that, unlike with MySQL & Postgre' SQL, You don't even need to install.

---

Here how SQLite3 to Connect to an SQLite database. in Python.

```
Python
1    import sqlite3
2    from sqlite3 import Error
3
4    def create_connection(path):
5        connection = None
6        try:
7            connection = sqlite3.connect(path)
8            print("Connection to SQLite DB Successful")
9        except Error as e:
10           print(f"The error '{e}' occured")
11
12       return connection
```

- Line 1 & 2 import sqlite3 & The Module Error Class
- Line 4 define function. create_connection()
- Line 7 uses .connect() from the sqlite3 module & take the SQLite database path a parameter
- Line 8 Print the Status of the Successful Connection.

---

- Line 9 Catches any exception that might be thrown if .connect() fails to establish a connection
- Line 10 displays the error message in the Console.
- The following script create a Connection to the SQLite DB.

```
connection = create_connection("E:\\sm-
app.sqlite")
```

### MySQL:

→ Unlike SQLite, there's no default Python SQL module that you can use to connect to a MySQL DB.

→ Instead, need to install a Python SQL driver for MySQL in Order to interact with a MySQL

→ Download this Python SQL Module with PiP

```
$ pip install mysql-connector-python
```

MySQL database has a two step process.

1. Make a Connection to a MySQL Server

2. Execute a Separate query to Create DB.

a function that Connects to MySQL database Server. & Return the Connection object.

**Python**

```
1.  import mysql.Connector
2.  from mysql.Connector import Error
3.
4.  def Create_Connection (host_name, user name,
                            user_Password):
5.      Connection = None
6.      try:
7.          Connection = mysql.Connector.Connect (
8.              host = host_name,
9.              user = user_name,
10.             Passwrd = user_Password
11.         )
12.         Print ("Connection to MySQL DB
                    Successful")
13.     except Error as e:
14.         Print (f" The error '{e}' occured)
15.
16.     return Connection
17.
18.  Connection = Create_Connection ("Localhost",
                    "root", " ")
```

From the above Script,
function Create_Connection() that
accept three parameter
    1. host_name
    2. user_name
    3. User-Password

---

**PostgreSQL:**

Like MySQL, there is no default Python SQL library that you can use to interact with a postgreSQL DB. Instead need to install a **3rd Party Python SQL driver** to interact with PostgreSQL.

One such Python SQL driver for PostgreSQL is **Psycopg2.**

Execute the following Command on your terminal to install the **Psycopg2 Python SQL module**

```
$. Pip install Psycopg2
```

Make a Connection with your PostgreSQL database

**Python**

```
1.  import Psycopg2
2.  from Psycopg2 import OperationalError
3.
4.  def Create_Connection (db_name, db_user,
                            db_host, db_Port):
5.      Connection = None
6.      try:
7.          Connection = Psycopg2.Connect (
                database = db_name,
                user = db_user
                Password = db_Password,
```

---

```
                host = db_host,
                Port = db_Port,
            )
8.
9.          Print ("Connection to PostgreSQL
                    DB Successful)
10.     . except OperationalError as e:
            Print (f" The error '{e}' occurred')
11.
12.     return Connection
```

**Creating Tables:**

In this Section, how to Create table inside there three databases. & Few table are
    1. Users
    2. Posts
    3. Comments
    4. Likes.

**SQLite:**

• To execute query in SQLite, use Cursor.execute().

• To define a function execute_query()

• execute() can execute any query passed to it

**Note:**
This Script should be executed in the Same file where you Created the Connection for your SQLite DB.

⑥

## Python

```
1.  def execute_query (Connection, query):
2.      Cursor = Connection. Cursor()
3.      try:
4.          Cursor. execute (query)
5.          Connection. Commit()
6.          Print (" Query execute Successfully")
7.      except Error as e:
8.          Print (f "The error "{e}" occured)
```

This Code execute the given query and Prints an error message

## Python

```
1.  Create_users_table = """
    CREATE TABLE IF NOT EXITS users (
    id INTEGER PRIMARY KEY
    AUTO_INCREMENT, name TEXT NOT NULL,
    age INTEGER, gender TEXT,
    nationality TEXT
    );
    """""
```

above Code Says to Create a table users
with following five Columns:

1. id    2. name    3. age
4. gender    5. nationality.

## MySQL:

Create tables in MySQL.
Just like with SQLite, need to pass your query to Cursor. execute(), Which is Returned by Calling .cursor() on the Connection Object.

## Python

```
1  def execute_query (Connection, query):
2      Cursor = Connection. query()
3      try:
4          Cursor. execute (query)
5          Connection. Commit()
6          Print (" Query execute Successful)
7      except Error as e:
8          Print (f "The Error '{e}" occures)
```

In line 4. Pass query to Cursor. execute().

## Python

```
1.  Create_users_table = """
2.  CREATE TABLE IF NOT EXITS users (
        id INT AUTO_INCREMENT,
        name TEXT NOT NULL,
        age INT, gender TEXT,
        nationality TEXT, PRIMARY KEY (id).
3.  ) ENGINE = InnoDB
    """
```

execute_query (Connection, Create_user_table)

MySQL uses the AUTO_INCREMET keyword to Create Column, the Values are automatically incremented.

## PostgreSQL:

→ Like with SQLite & MySQL database, the Connection object returned by Psycopg2. Connect() cursor object.

→ To execute Python-SQL queries on your PostgreSQL DB.

Define a function execute_query():

## Python

```
1  def execute_query (Connection, query):
2      Connection. autoCommit = True
3      Cursor = Connection. Cursor()
4      try:
5          Cursor. execute (query)
6          Print ("Query Executed Successfully")
7      except OperationalError as e:
8          Print(f" The error "{e} occured)
```

## Python

```
1.  Create_user_table = """
2.  CREATE TABLE IF NOT EXITS User (
    id SERIAL PRIMARY KEY
    name TEXT NOT NULL, age INTEGER,
    gender TEXT, nationality TEXT
    ) """
```

execute_query (Connection, Create_user_table)

users Table PostgreSQL is slightly different than SQLite & MySQL.
The keyword SERIAL is used to Create Column that Increment automatically

## Updating Table Records:

In this Section covers the Process for updating records using Python SQL libraries for SQLite, PostgreSQL, & MySQL

### SQLite:

As an example, you can update the description of the post with an id # 2

#### Python

```
1   Select_Post_description = "SELECT description
                               From POST WHERE id=2"
2   Post_description = execute_read_query(Connection,
                       Select_Post_description)
3   for description In Post_description:
4       Print (description)
```

#### Python

```
1.  update_Post_description = " "
2       UPDATE
3           POSTS
4       SET
            description = "The Weather has become
                          Pleasant now"
5       WHERE
6           id = 2  ""
execute_query (Connection, update_Post_description)
```

execute the query again, & see the following Result shell

("The Weather has become Pleasant now",)

### MySQL:

The Process of updating records in MySQL with mysql-Connector-Python is also a carbon of the Sqlite3 PythonSQL Module.

For example, the following Script updates the description of the Post with an id#2

#### Python

```
1   update_Post_description = " "
2       UPDATE
3           posts
4       SET
            description = "The Weather has become
                           Pleasant now"
5       WHERE
6           id = 2
    " "
execute_query (Connection, update_post_description)
```

### PostgreSQL:

The update query for PostgreSQL is Similar to what you have Seen with SQLite and MySQL. Use above Scripts to update record in your PostgreSQL table.

## Deleting Table Records:

To delete the table records using PythonSQL modules for SQLite, MySQL & PostgreSQL. The Process of deleting records is uniform for all three database

### SQLite:

To delete records from your SQLite database, execute_query() will Create a cursor object using the Connection & pass the String to Cursor.execute(), which will delete the records.

#### Python

```
1   delete_Comment = "DELETE FROM Comments
                      WHERE id=5"
2   execute_query (Connection, delete_Comment)
```

### MySQL:

The Process of deletion in MySQL is also similar to SQLite. Ex:

#### Python

```
1   delete_Comment = "DELETE FROM Comments
                      WHERE id=2"
2   execute_query (Connection, delete_Comment)
```

### PostgreSQL:

The delete query for PostgreSQL is also Similar to SQLite & MySQL. Write a delete query String by using the DELETE keyword and then passing the query & the Connection object to execute_query(). This will delete the specified records from your PostgreSQL DB.

## Conclusion:

Three Common Python SQL libraries, Sqlite3, mysql-Connector-Python & Psycopg2 allow you to connect a Python application to SQLite, MySQL & PostgreSQL

- Interact
- Use
- Execute

## SQLite

### Supported Data Types

- null
- integer
- real
- text
- blob

## MySQL

### Supported Data Types

- tinyint
- smallint
- mediumint
- int (or) integer
- bigint
- float
- dec, decimal, fixed (or) numeric
- bool or boolean
- bit

### Date and Time Types

- date
- datetime
- timestamp
- time
- year

### String Types

- char
- varchar
- binary
- varbinary
- blob
- tiny blob
- medium blob
- long blob
- text
- tiny text
- medium text
- long text
- enum
- set

## PostgreSQL

### Numeric Types

- bigint
- bigserial
- double precision
- integer
- numeric (or) decimal
- real
- smallint
- small serial
- serial

### Character Types

- character
- character varying (or) varchar
- text

### Date and Time Types

- date
- interval
- time (or) time without time zone
- time with time zone
- timestamp (or) timestamp without time zone
- timestamp with time zone

### Geometric Types

- bot
- circle
- line
- lseg
- path
- point
- polygon

### Network address Types:
- Cidr
- inet
- macaddr

### Bit String Types:
- bit
- bit varying

### Text Search Types:
- tsquery
- tsvector

### JSON Types:
- json
- jsonb

### Other Data Types:
- boolean
- bytea
- money
- pg_lsn
- uuid
- xml
- txid-snapshot

↳ A user level transaction ID Snapshot

9

# Unit III — Data Cleanup

## Why Clean Data?

↳ It is the process of fixing or removing incorrect, Corrupted, incorrectly formatted, duplicate or incomplete data within a dataset.

## Basics

1. Remove duplicate or irrelevant observations
2. Fix Structural errors
3. Filter unwanted outliers
4. Handle missing data
5. Validate and QA

- Substitute missing data, such as N/A for null categories or O for numerical values.
- Substitute average numbers for the missing numerical values.

## Identifying Values for Data Clean up.

* Remove duplicates
* Remove irrelevant data
* Standardize capitalization
* Convert data type
* Clear formatting
* Fix errors
* Language Translation
* Handle missing values

## Data Formatting

↳ Converting data to the required format

↳ When merging data from several sources, check all the variables inside a particular attribute are written consistently.

↳ Formats for dates, money (4.03, $4.03, or even $4.03),

addresses, and so on.

↳ Across the dataset, the input format should be constant.

## Finding Outliers and Bad data

1. Sort your data
2. Graph your data
3. Calculate the Z-Score

$$Z = (x - \mu) \div \sigma$$

where

$x$ = Raw measurement

$\mu$ = the mean

$\sigma$ = the standard deviation

4. Calculate the interquartile range

$$IQR = Q_3 - Q_1$$

$Q_3$ = the third quartile = the median of the upper half of the dataset.

$Q_1$ = the first quartile = the median of the lower half of the dataset.

High outlier $\geq Q_3 + (1.5 \times IQR)$

Low outlier $\leq Q_1 - (1.5 \times IQR)$

5. Use a hypothesis test.

## Finding Duplicates :-

↳ The process of duplicate detection is preceded by a data preparation stage, during which data entries are stored in an uniform manner in the database.

## Fuzzy Matching

↳ identifies similar, but not identical elements in data table sets.

## Levenshtein distance :-

1. Power → Powder (insertion of "w")
   — Distance := 1

2. Lovin → Loving (insertion of "g")
   — Distance := 1

3. Porpoise → Purpose (insertion of "i", substitution "o" for "u") — Distance := 2

## Hamming Distance :-

Eg :- "Corn" and "Cork"

The binary ASCII character table assigns the code (01101110) for N and (01101011) for K, the difference between each letter's code occurs in two locations, thus an HD of 2.

## Damerau - Levenshtein

String 1 : Micheal

String 2 : Michaela

Operation 1 :- Transposition : swap characters "a" and "e"

Operation 2 :- insert "a" (end of string 2)

Distance = 2

Each operation has a count of "1", so each insertion, deletion, transposition, etc is weighted equally.

## Regular Expressions (Regex)

↳ is extremely and amazingly powerful is searching and manipulating text strings, particularly is processing text files.

* character
* Special Regex characters
* Escape Sequences
* A sequence of characters
* OR operator (1)
* character class
* Occurrence Indicators (or Repetition operators)
* Metacharacters
* Position Anchors
* Parenthesized Back References
* Laziness (Curb Greediness for Repetition operators) :
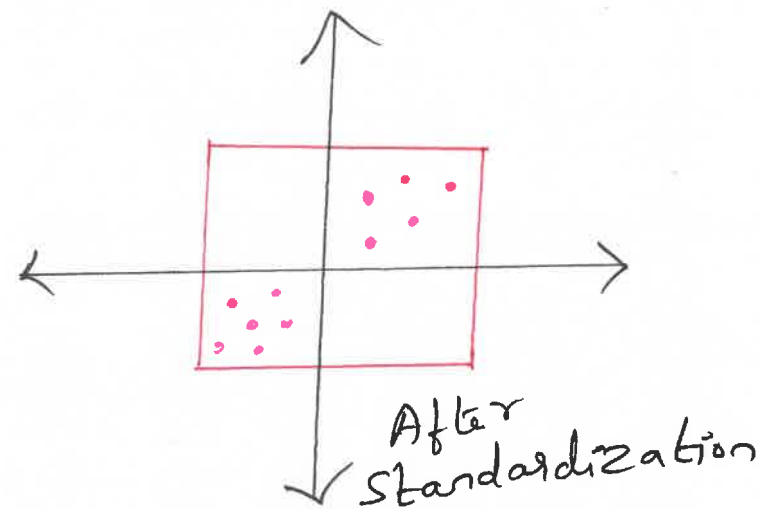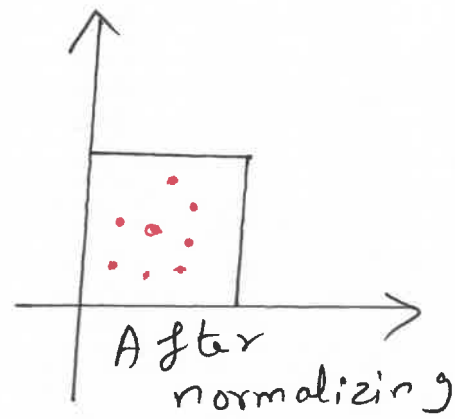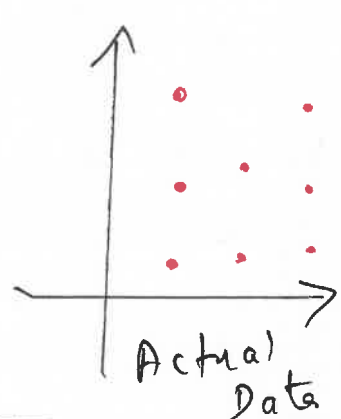* ? , + ? , ?? , {m, n} ? (m, 3 ?

# Normalization and Standardization

- is a part of data processing and cleansing techniques.
- is to make the data homogenous overall records and fields.
- helps in creating a linkage between the entry data which in turn helps in cleaning and improving data quality.

## Standardization is the process of placing dissimilar features on the same scale.

- Rescaling the attributes is such a way that their mean is 0 and standard deviation becomes 1.



Actual Data

After normalizing



After Standardization

## Saving Data:

Pickle module, which is a part of the standard library to save data in Python.

```
import pickle
Class MyClass ():
    def __init__(self, param):
        self.param = param
def save_object (obj):
    try:
        with open ("data.pickle",
                    "wb") as f:
            pickle.dump (obj, f,
```

Protocol = Pickle.HIGHEST_
                        PROTOCOL)

```
    except Exception as ex:
        print(" Error during
            pickling object (Possibly
            unsupported):", ex)
obj = MyClass (10)
save_object (obj)
```

If you run this script, a file called data.pickle, which contains the saved data.

## Determining Suitable Data cleanup

1. Look into your data

2. Look at the proportion of missing data

3. Check the data type of each column

12

4. If columns of strings are available, check for trailing whitespaces.

5. Dealing with Missing values (NaN values)

6. Extracting more information from your dataset to get more variables.

7. Check the unique values of columns.

## Scripting the cleanup and Testing with New Data.

The following script is both a testing and cleaning script to test for all lowercase letters, and to lowercase them if they are not already lowercase.

```
aimReturnStatus = True;
aimReturnValue = aimValue;
if (aimValue != null)
{
lower = aimValue.toLowerCase();
if (! aimValue.equals(lower))
{
aimReturnStatus = false;
aimReturnValue = lower;
}
}
```

When applied as a test script, the following script flags longer than 16 characters as failed.

When applied as a cleaning script, it truncates the string and replaces the last 3 characters with "...".

```
limit = 16;

aimReturnStatus = True;
aimReturnValue = aimValue;
if (aimValue != null)
{
len = aimValue.length;
if (len > limit)
{
aimReturnStatus = false;
aimReturnValue = aimValue.substring(0, limit-3)
+ "...";
}
}
```

Data cleaning is the process by which you correct errors and anomalies in the case data.

↳ Save your time from having to do a code review and make corrections and also from having to enter a commit message each time

# UNIT IV

## Exploring Data:-
↳ is the first step of data analysis.
↳ used to explore and visualize data

## Importing Data
- Delimited text files (.txt)
- Comma Separated values text files (.csv)

## Exploring Table Functions
- Head
- Info
- Describe
- Shape
- Columns & Index
- astype

## Data Exploration and Analysis.

- Sort Values
- Value - Counts
- drop - duplicates
- Subsetting
- Adding a new Column
- Set Index and Reset index
- loc and iloc
- groupBy
- Pivot Table

## Joining numerous datasets:-
The three best and most time saving ways to combine multiple datasets using Python pandas methods.

- merge() - To combine the datasets on common column or index or both.
- concat() - To combine the datasets across rows or columns
- join() - To combine the datasets on key column or index

## Identifying Correlations
The statistics. correlation() method in python is used to return Pearson's correlation coefficient between two inputs.

$$r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

where

$r_{xy}$ = the correlation coefficient of the linear relationship between the variables x and y.

$x_i$ = the values of the x-variable in a sample

$\bar{x}$ = the mean of the values of the x-variable.

$y_i$ = the values of the y-variable in a sample.

$\bar{y}$ = the mean of the values of the y-variable

14

# Methods to detect Outliers in Python.

* Z-Score
* Interquartile Range (IQR)
* Tukey's fences.

## Creating Groupings

The groupby() function in pandas is used to group DataFrame rows based on one or more Columns.

* Import the required libraries
* Create a sample dataframe
* Group the dataframe based on a Condition
* Group the dataframe based on multiple Conditions

# Data Analysis

↳ is the technique of collecting, transforming, and organizing data to make future predictions and informed data-driven decision.

* Specify Data Requirement
* Prepare or Collect Data
* Clean and Process
* Analyze
* Share
* Report

## Dataset Splitting:

sklearn is the most useful and Robust library for machine learning in python.

### Parameters:-

* arrays
* test-size
* train-size
* Random-state
* shuffle
* stratify

**Arrays:-** inputs such as lists, arrays, data frames or matrices

**test-size:-** is a float value whose value ranges between 0.0 and 1.0

**train-size:-** is a float value whose value ranges between 0.0 and 1.0

**Random-state:-** Controls the shuffling applied to the data before applying the split

**Shuffle:-** is used to shuffle the data before splitting.

**Stratify:-** split the data is a stratified fashion

15

# Data Visualization in Python

* Line plots
* Box plots
* Heatmaps
* Lag Plots
* AutoCorrelation plots.

Analytics Software is software with a source code that anyone can inspect, modify or enhance.

Python libraries for Data Visualisation:

* Matplotlib
* Plotly
* Seaborn
* GGplot
* Altair

* Bokeh
* Pygal
* Geoplotlib
* Folium
* Gleam

## Publishing the Data-open Source Platform.

The steps to publish a Python package are follows:-

* Write your python module and save it under a directory
* Create the setup.py file with necessary information
* Choose a LICENSE and provide a README file for your project
* Generate the distribution archives on local machine
* Try installing the package on your local machine.

* Publish the package to the TestPyPi repository to check if all works well.
* Finally, publish the package to the PyPi repository.

## Open-Source Platform

* Seal Report
* KNIME
* Knowage
* Metabase
* Microsoft Power BI Desktop
* Rapid Miner
* Report Server
* Tableau
* Web Data Rocks
* Zoho Analytics
* Dataiku DSS Free Edition

16

# UNIT V  VISUALIZING DATA WITH PANDAS AND MATPLOTLIB

## An Introduction to Matplotlib

↳ is a low level graph plotting library in python that serves as a visualization utility.

↳ Created by John D. Hunter

↳ is an open source and we can use it freely.

↳ mostly written in python, few segments are in C, and Javascript for platform compatibility

### Basics:

↳ used to create 2D graphs and plots by using python script.

↳ uses **pyplot** which provides feature to control line styles, font properties, formatting axes, etc.

↳ supports a wide variety of graphs and plots namely - histogram, bar charts, power spectra, error charts etc.
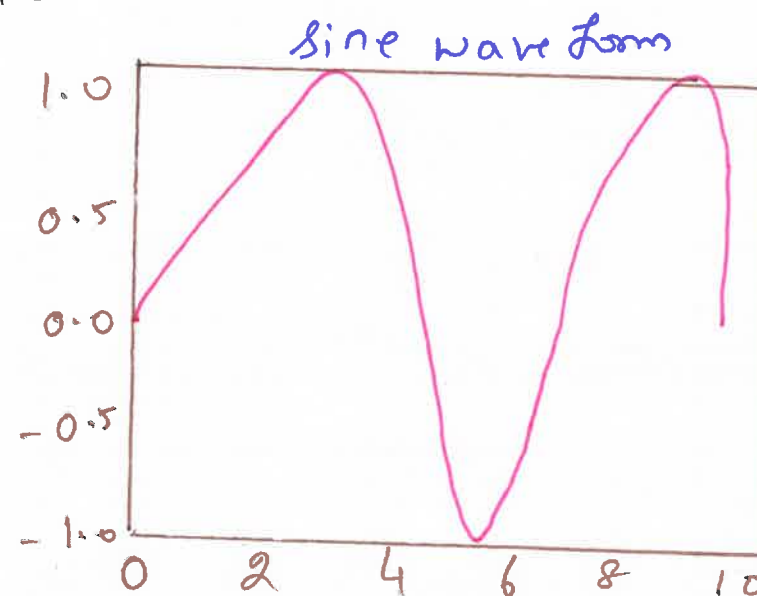
### Example:

```
import numpy as np
import matplotlib.pyplot as plt
x = np.arange( 0, 3 * np.pi, 0.1)
y = np.sin (x)
plt.title (" sine wave form")
plt.plot (x, y)
plt.show ()
```



Sine wave form

## Plot Components

axis - Axes of the graph that are typically x, y and z (for 3D plots)

title - Title of the plot.

label - Name for the axis (e.g x label or y label).

legend - Legend for the plot.

ticklabel - Text or Values represented on the axis.

Symbol - Symbol for data points that can be presented with different symbol shapes/colors.

Size - Size of, for example, a point on a scatter plot.

linestyle – The style how the line should be drawn. Can be solid or dashed, for example

line width – The width of a line on a plot.

alpha – Transparency level of a filled element in a plot.

tick(s) – Refers to the tick marks on a plot.

annotation – Refers to the text added to a plot.

padding – The distance between an (axis/tick) label and the axis.

## Plotting with Pandas

| Parameter | Type |
|-----------|------|
| Kind | String |
| x/y | String (or list |
| ax | Axes |

| | | | |
|---|---|---|---|
| Subplots | Boolean | xticks/yticks | List of values |
| layout | Tuple of (rows, columns) | rot | Integer |
| figsize | Tuple of (width, height) | xlim/ylim | Tuple of the form (min, max) |
| title | string for the plot title or a list of strings for subplot titles | sharex/sharey | Boolean |
| legend | Boolean | fontsize | Integer |
| label | String if a single column is being plotted. | grid | Boolean |
| style | string if a single column is being plotted. | | |
| Color | String if a single Column is being plotted. | | |
| Colormap | String or matplotlib colormap object | | |
| logx/logy/log log | Boolean | | |

## Plotting

↳ Pandas uses the plot() method to create diagrams.

## Scatter Plot

↳ Specify that you want a scatter plot with the kind argument.
Kind = 'scatter'
A scatter plot needs an x and y axis.

18

# Lag Plot

↳ checks whether a dataset or time series is Random or not.

↳ Random data should not exhibit any identifiable structure in the lag plot.

↳ Non Random structure in the lag plot indicates that the underlying data are not Random.

## Auto Correlation Plots :

↳ is a Commonly used tool for checking Randomness in a data set.

## Syntax :

matplotlib.pyplot.accor(x, *, data = None, ** kwargs)

## Parameters :

x, detrend, normed, userlines, maxlags, linestyle and marker.

# Bootstrap Plot :-

↳ is used to eliminate the uncertainty of a statistic by relying on Random sampling with replacement.

↳ This function will generate bootstrapping plots for Mean, Median and mid-range statistics for the given number of samples of the given size.

## Counts and Frequencies

↳ Counting the frequency of occurrence of a word in a body of text is often needed during text processing.

import nltk
from nltk. Corpus import brown

CFd = nltk. Conditional Freq Dist ( (genre, word)
for genre in brown. Categories ()
for word in brown. words( Categories = genre))
Categories = [ 'hobbies', 'romance', 'humor']
Searchwords = [ 'may', 'might', 'must', 'will']

CFd. tabulate (Conditions = Categories, samples = Searchwords)

## Output:

| | may | might | must | will |
|---|---|---|---|---|
| hobbies | 131 | 22 | 83 | 264 |
| romance | 11 | 51 | 45 | 43 |
| humor | 8 | 8 | 9 | 13 |

# ① Analyzing Road Safety in the UK

* Evaluate the median severity value of accident caused by various Motorcycles.

* Evaluate Accident Severity and Total Accidents per vehicle type.

* Calculate the average severity by vehicle type

* Calculate the average severity and Total Accidents by Motor Cycle.

# ② International Debt Statistical Analysis

* What is the total amount of money that Countries owe to the World Bank?

* Which Country has the highest debt, and how much is that?

* What is the mean debt owed by Countries for different debt indicators?

# ③ Analyzing the World Population

* Which Country has the highest population?

* Which Country has the least number of people?

* Which Country is witnessing the highest population growth?

* Which Country has an extraordinary number for the population?

* Which is the most densely Populated Country in the world?

# ④ Analyzing Sales Data

* What was the total revenue generated by the Retailer during the time period covered by the dataset?

* Which product had the highest total revenue?

* Which Customers generated the most revenue for the retailer?

* How many unique product were sold?

# ⑤ Inventory Control Management

* Which year had the highest sales?

* How was the weather during the year of highest sales?

* Conclude whether the weather has an essential impact on sales.

* Do the sales always rise nearthe holiday season for all the years?

* Analyze the relationship between sales and the different macroeconomic variables in the dataset.