# TRAINING REPORT

## At

## SATHYABAMA INSTITUTE OF SCIENCE AND TECHNOLOGY

## (DEEMED TO BE UNIVERSITY)

Submitted in partial fulfillment of the requirements for the award
of Bachelor of Engineering Degree in

Computer Science and Engineering

By

## Kumar Aditya (38110276)



## DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

## SCHOOL OF COMPUTING

## SATHYABAMA INSTITUTE OF SCIENCE AND TECHNOLOGY

## JEPPIAAR NAGAR, RAJIV GANDHI SALAI,

## CHENNAI – 600119, TAMILNADU

## AUGUST 2020

SATHYABAMA

INSTITUTE OF SCIENCE AND TECHNOLOGY
(DEEMED TO BE UNIVERSITY)
**Accredited with Grade "A" by NAAC**
(Established under Section 3 of UGC Act, 1956)
**JEPPIAAR NAGAR, RAJIV GANDHI SALAI, CHENNAI– 600119**
www.sathyabamauniversity.ac.in

## DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

### BONAFIDE CERTIFICATE

This is to certify that this Project Report is the bonafide work of **Kumar Aditya (38110276)** who carried out the project entitled "**Covid – 19 Data Processing**" under my supervision from April 2020 to July 2020.

**Internal Guide**

**Dr. Albert Mayan J, M.E., Ph.D.**

**Head of Department**

**Dr. S.VIGNESHWARI, M.E., Ph.D.,**
**Dr. Lakshmanan L, M.E, Ph.D.,**

Submitted for Viva voce Examination held on_____

**Internal Examiner**                                            **External Examiner**

# DECLARATION

I, **Kumar Aditya (38110276)** hereby declare that the Project Report entitled "**Covid – 19 Data Processing**" done by me under the guidance of **Dr. Albert Mayan J, M.E., PhD** at Sathyabama Institute of Science and Technology is submitted in partial fulfillment of the requirements for the award of Bachelor of Engineering degree in Computer Science and Engineering.

**DATE:**

**PLACE: CHENNAI**                                                   **SIGNATURE OF THE CANDIDATE**

# ACKNOWLEDGEMENT

I am pleased to acknowledge my sincere thanks to **Board of Management** of **SATHYABAMA** for their kind encouragement in doing this project and for completing it successfully. I am grateful to them.

I convey my thanks to **Dr. T.Sasikala M.E., Ph.D.**, **Dean**, School of Computing , **Dr.S.Vigneshwari M.E., Ph.D., and Dr.L.Lakshmanan M.E., Ph.D.,** Heads of the Department of Computer Science and Engineering for providing me necessary support and details at the right time during the progressive reviews.

I would like to express my sincere and deep sense of gratitude to my Project Guide **Dr. Albert Mayan, M.E., PhD** for his valuable guidance, suggestions and constant encouragement paved way for the successful completion of my project work.

 I wish to express my thanks to all Teaching and Non-teaching staff members of the **Department of Computer Science and Engineering** who were helpful in many ways for the completion of the project.

# TRAINING CERTIFICATE

**coursera**

**4 Courses**

What is Data Science?

Tools for Data Science

Data Science Methodology

Databases and SQL for Data Science
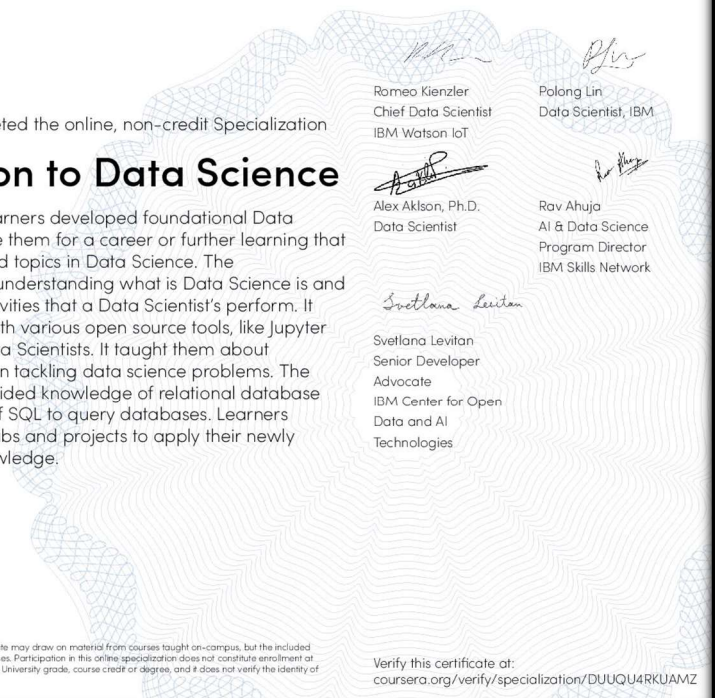
**IBM.**

01/10/2020

**Kumar Aditya**

has successfully completed the online, non-credit Specialization

## Introduction to Data Science

In this Specialization learners developed foundational Data Science skills to prepare them for a career or further learning that involves more advanced topics in Data Science. The specialization entailed understanding what is Data Science is and the various kinds of activities that a Data Scientist's perform. It familiarized learners with various open source tools, like Jupyter notebooks, used by Data Scientists. It taught them about methodology involved in tackling data science problems. The specialization also provided knowledge of relational database concepts and the use of SQL to query databases. Learners completed hands-on labs and projects to apply their newly acquired skills and knowledge.

Romeo Kienzler
Chief Data Scientist
IBM Watson IoT

Polong Lin
Data Scientist, IBM

Alex Aklson, Ph.D.
Data Scientist

Rav Ahuja
AI & Data Science
Program Director
IBM Skills Network

Svetlana Levitan
Senior Developer
Advocate
IBM Center for Open
Data and AI
Technologies

Verify this certificate at:
coursera.org/verify/specialization/DUUQU4RKUAMZ

4

# ABSTRACT

Covid – 19 Data Processing is a Program which can load, merge, clean and aggregate the COVID-19 time series data. This program can give graphs and stats about any country in the world. It is designed by using Python Language, Matplotlib, Seaborn, Alter, Numpy and Pandas.

It takes the updated data from the web and show the current stats of Covid – 19 such as Confirmed, Recovered, Death and Active cases around each country for which the user requires. Further this program also compares data of newly confirmed cases between multiple countries.

# TABLE OF CONTENTS

# LIST OF FIGURES

# CHAPTER 1

## 1.1 OUTLINE

This Covid – 19 Data Processing which is made using Python Language, Matplotlib, Seaborn, Alter, Numpy and Pandas can be used for creating graphs and stats for countries and also compare data of countries. Pandas is used to create DataFrame. Matplotlib, Seaborn and Alter is used to create graphs and stats table.

## 1.2  WHY USE PYTHON, NUMPY, PANDAS

Python is open source, interpreted, high level language and provides great approach for object-oriented programming. It is one of the best languages used by data scientist for various data science projects/application. It provides great functionality to deal with mathematics, statistics and scientific function. It provides great libraries to deals with data science application. It uses the elegant syntax; hence the programs are easier to read. The interactive mode of Python makes its simple to test codes. It allows developer to run the code anywhere, including Windows, Mac OS X, UNIX, and Linux. It is free software in a couple of categories. It does not cost anything to use or download Pythons or to add it to the application.

Numpy is Python library that provides mathematical function to handle large dimension array. It provides various method/function for Array, Metrics, and linear algebra. NumPy stands for Numerical Python. It provides lots of useful features for operations on n-arrays and matrices in Python. The library provides vectorization of mathematical operations on the NumPy array type, which enhance performance and speeds up the execution. It's very easy to work with large multidimensional arrays and matrices using NumPy. Nearly every scientist working in Python draws on the power of NumPy. NumPy brings the computational power of languages like C and Fortran to Python, a language much easier to learn and use. With this power comes simplicity: a solution in NumPy is often clear and elegant. The core functionality of NumPy is its "ndarray", for *n*-dimensional array, data structure. These arrays are strided views on memory.

Pandas is a fast, powerful, flexible and easy to use open source data analysis and

manipulation tool, built on top of the Python programming language. Pandas is one of the most popular Python library for data manipulation and analysis. Pandas provide useful functions to manipulate large amount of structured data. Pandas provide easiest method to perform analysis. It provides large data structures and manipulating numerical tables and time series data. Pandas is a perfect tool for data wrangling. Pandas is designed for quick and easy data manipulation, aggregation, and visualization. There two data structures in Pandas: Series handle and store data in one-dimensional data whereas DataFrame handle and store Two-dimensional data. It is a group by engine allowing split-apply-combine operations on data sets.

Data science is an interdisciplinary field focused on extracting knowledge from data sets, which are typically large (see big data). The field encompasses analysis, preparing data for analysis, and presenting findings to inform high-level decisions in an organization. As such, it incorporates skills from computer science, mathematics, statistics, information visualization, graphic design, complex systems, communication and business. Data science is a "concept to unify statistics, data analysis, machine learning, domain knowledge and their related methods" in order to "understand and analyze actual phenomena" with data.

## 1.3 WHY USE MATPLOTLIB, SEABORN, ALTAIR

Data science is a "concept to unify statistics, data analysis. Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python. It is useful Python library for Data Visualization. Descriptive analysis and visualizing data is very important for any organization. It provides various method to Visualize data in more effective way. It is a comprehensive library for creating static, animated, and interactive visualizations in Python. One of the greatest benefits of visualization is that it allows us visual access to huge amounts of data in easily digestible visuals. It consists of several plots like line, bar, scatter, histogram etc.

Seaborn is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics. It is a library for making statistical graphics in Python. It is built on top of matplotlib and closely integrated with pandas data structures. It is a Specialized support for using categorical variables to

show observations or aggregate statistics. It has automatic estimation and plotting of linear regression models for different kinds dependent variables. It has high-level abstractions for structuring multi-plot grids that let you easily build complex visualizations It has tools for choosing color palettes that faithfully reveal patterns in your data. Its specialized support for using categorical variables to show observations or aggregate statistics.

Altair is a declarative statistical visualization library for Python. With Altair, you can spend more time understanding your data and its meaning. Altair's API is simple, friendly and consistent and built on top of the powerful Vega-Lite visualization grammar. This elegant simplicity produces beautiful and effective visualizations with a minimal amount of code. Altair offers a powerful and concise visualization grammar that enables you to build a wide range of statistical visualizations quickly. It can be used to change several other attributes of the tables such as the name of the table and name of the columns. It can changes the structure or schema of an already created database table by adding or removing one or more columns.

## 1.4 LITERATURE REVIEW

Python was conceived in the late 1980s by Guido van Rossum at Centrum Wiskunde & Informatica (CWI) in the Netherlands as a successor to the ABC language (itself inspired by SETL), capable of exception handling and interfacing with the Amoeba operating system. Its implementation began in December 1989. Van Rossum shouldered sole responsibility for the project, as the lead developer, until 12 July 2018, when he announced his "permanent vacation" from his responsibilities as Python's *Benevolent Dictator For Life*, a title the Python community bestowed upon him to reflect his long-term commitment as the project's chief decision-maker. He now shares his leadership as a member of a five-person steering council. In January 2019, active Python core developers elected Brett Cannon, Nick Coghlan, Barry Warsaw, Carol Willing and Van Rossum to a five-member "Steering Council" to lead the project.

It is a multi-paradigm programming language. Object-oriented programming and structured programming are fully supported, and many of its features support functional programming and aspect-oriented programming (including by meta programming and

metaobjects (magic methods)). Many other paradigms are supported via extensions, including design by contract and logic programming.

Python uses dynamic typing and a combination of reference counting and a cycle-detecting garbage collector for memory management. It also features dynamic name resolution (late binding), which binds method and variable names during program execution.

Python's design offers some support for functional programming in the Lisp tradition. It has filter, map, and reduce functions; list comprehensions, dictionaries, sets, and generator expressions. The standard library has two modules (itertools and functools) that implement functional tools borrowed from Haskell and Standard ML and Data Science.

In 1962, John Tukey described a field he called "data analysis," which resembles modern data science. Later, attendees at a 1992 statistics symposium at the University of Montpellier II acknowledged the emergence of a new discipline focused on data of various origins and forms, combining established concepts and principles of statistics and data analysis with computing.

The term "data science" has been traced back to 1974, when Peter Naur proposed it as an alternative name for computer science. In 1996, the International Federation of Classification Societies became the first conference to specifically feature data science as a topic. However, the definition was still in flux. In 1997, C.F. Jeff Wu suggested that statistics should be renamed data science. He reasoned that a new name would help statistics shed inaccurate stereotypes, such as being synonymous with accounting, or limited to describing data. In 1998, Chikio Hayashi argued for data science as a new, interdisciplinary concept, with three aspects: data design, collection, and analysis.

During the 1990s, popular terms for the process of finding patterns in datasets (which were increasingly large) included "knowledge discovery" and "data mining. The modern conception of data science as an independent discipline is sometimes attributed to William S. Cleveland. In a 2001 paper, he advocated an expansion of statistics beyond theory into technical areas; because this would significantly change the field, it warranted

a new name. Data science is a growing field. A career as a data scientist is ranked at the third best job in America for 2020 by Glassdoor, and was ranked the number one best job from 2016-2019.

As such, it incorporates skills from computer science, mathematics, statistics, information visualization, graphic design, complex systems, communication and business. Statistician Nathan Yau, drawing on Ben Fry, also links data science to human-computer interaction: users should be able to intuitively control and explore data. In 2015, the American Statistical Association identified database management, statistics and machine learning, and distributed and parallel systems as the three emerging foundational professional communities.

## 1.5 PROBLEM STATEMENT

To process data of Covid – 19 and plot graphs and stats using Data Visualization Tools. Such that one can get data of different countries of world and plot its graphs and stats like Confirmed, Recovered, Deaths etc. Also compare data of different countries.

## 1.6 OBJECTIVES

A DataFrame to create and process Covid – 19 data and give different graphs and stats of different countries with respect to Confirmed, Recovered, Death.

# CHAPTER 2

## AIM & SCOPE OF DATA SCIENCE

### 2.1 REQUIREMENTS

#### *2.1.1 HARDWARE REQUIREMENTS*

- Inter core i5 7th Gen / Ryzen 5 2500u
- 8GB of DDR4 RAM
- 128 GB SSD / HDD
- 2GB Nvidia GeForce RTX 1050 Ti

#### *2.1.2 SOFTWARE REQUIREMENTS:*

- Python installed on your system
- Jupyter Notebook

The Jupyter Notebook is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations and narrative text. Uses include: data cleaning and transformation, numerical simulation, statistical modeling, data visualization, machine learning, and much more.

### 2.2 ROLE OF JUPYTER

Jupyter is a free, open-source, interactive web tool known as a computational notebook, which researchers can use to combine software code, computational output, explanatory text and multimedia resources in a single document. Computational notebooks have been around for decades, but Jupyter in particular has exploded in popularity over the past couple of years. This rapid uptake has been aided by an enthusiastic community of user–developers and a redesigned architecture that allows the notebook to speak dozens of programming languages — a fact reflected in its name, which was inspired, according to co-founder Fernando Pérez, by the programming languages Julia (Ju), Python (Py) and R. Computational notebooks are essentially laboratory notebooks for

scientific computing.

For data scientists, that format can drive exploration. Notebooks, Barba says, are a form of interactive computing, an environment in which users execute code, see what happens, modify and repeat in a kind of iterative conversation between researcher and data. They aren't the only forum for such conversations — IPython, the interactive Python interpreter on which Jupyter's predecessor, IPython Notebook, was built, is another. But notebooks allow users to document those conversations, building "more powerful connections between topics, theories, data and results".

The Jupyter notebook has two components. Users input programming code or text in rectangular cells in a front-end web page. The browser then passes that code to a back-end 'kernel', which runs the code and returns the results.

Jupyter's newest variant is JupyterLab, which launched as a beta in January 2018 and is available (like the Jupyter notebook) either as a stand-alone package or as part of the free Anaconda scientific-computing environment.

## 2.2.1 LANGUAGES USED FOR PROJECT

Machine Learning , data science and artificial Intelligence-based projects are obviously what the future holds . We want better personalization, smarter recommendations, and improved search functionality. Our apps can see, hear, and respond-that's what artificial intelligence (AI) has brought , enhancing the user experience and creating value across many industries.

**FEATURES OF PYTHON:**

Following are some of the important features that make Python the first choice of Data Science Engineers.

**Simple and consistent:**
Python offers concise and readable code. While complex algorithms and versatile workflows stand behind data science, machine learning and AI, Python's simplicity allows

developers to write reliable systems. Python code is understandable by humans, which makes it easier to build models for machine learning.

**Extensive selection of libraries and frameworks:**
Implementing Data Science and ML algorithms can be tricky and requires a lot of time. It's vital to have a well-structured and well tested environment to enable developers to come up with the best coding solutions.
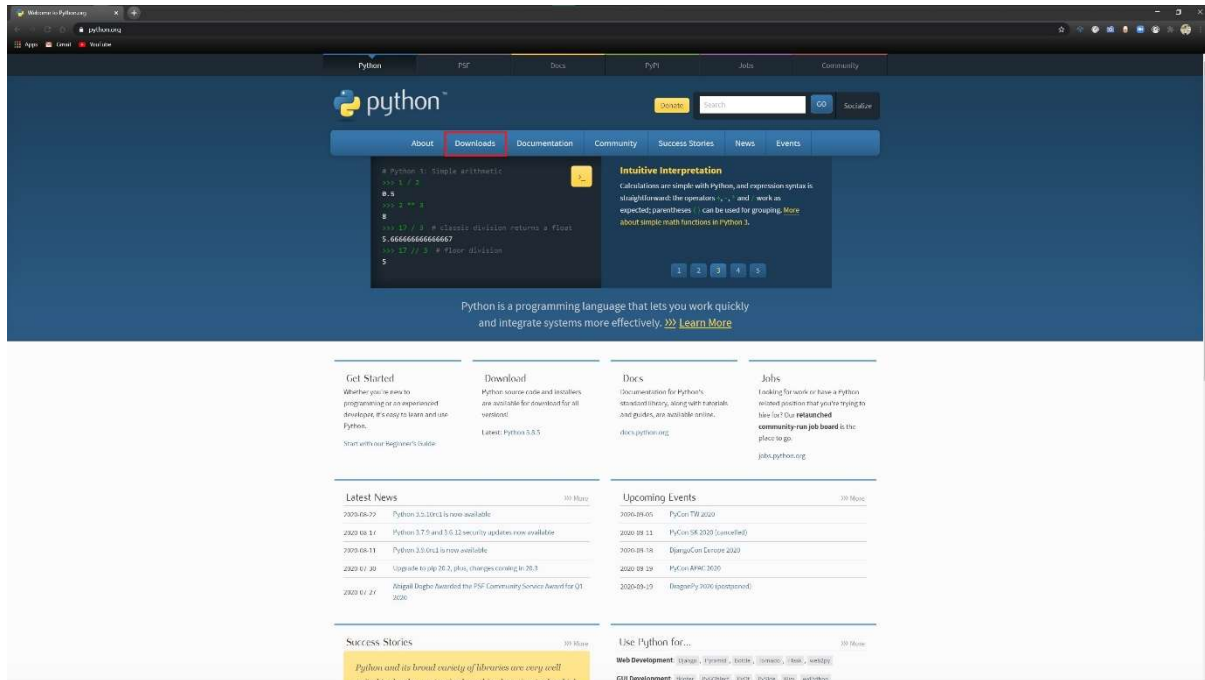
- Numpy and Pandas for general purpose data analysis.
- Matplotlib for creating static, animated, and interactive visualizations.
- Seaborn for high-level interface for drawing attractive and informative statistical graphics.
- Altair for understanding data and its meaning.

**Platform independence:**

Python software can be easily distributed and used on that operating system without a Python interpreter. Great Community and Popularity: In the Developer survey 2018 by Stack Overflow, Python was among the most popular programming language, which ultimately mean that you can find easily a development company.
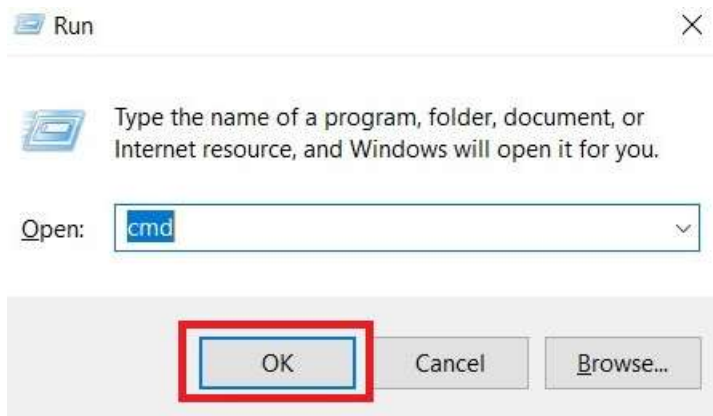
## 2.3 HOW TO INSTALL PYTHON

From Website:



You can download Python from their website directly to your system for development purposes.
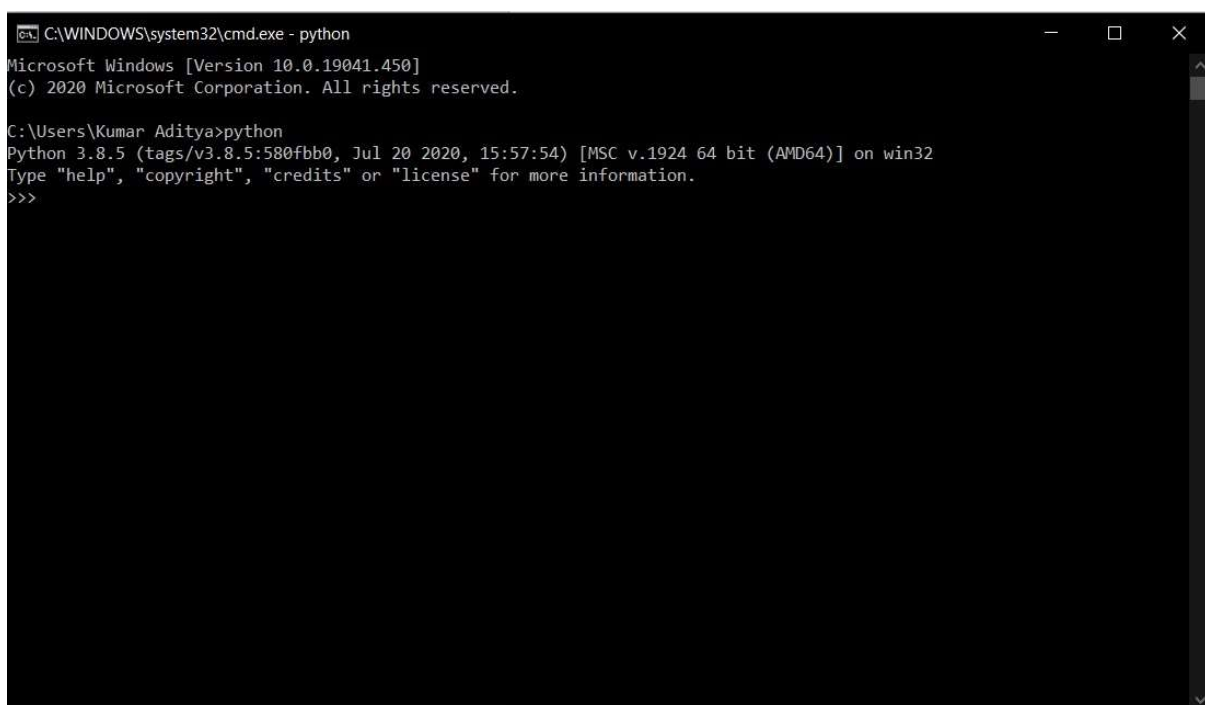


Do not forget to "Add Python to PATH"

❖ After installation is complete you will be able to see Python in Command Prompt.
   *Run -> cmd*



❖ You can see information about your version and switch to python by typing "python" in Command Prompt.

## 2.4 HOW TO INSTALL JUPYTER NOTEBOOK

Using pip:

If you use pip, you can install it with: (In Command Prompt)

pip install jupyterlab


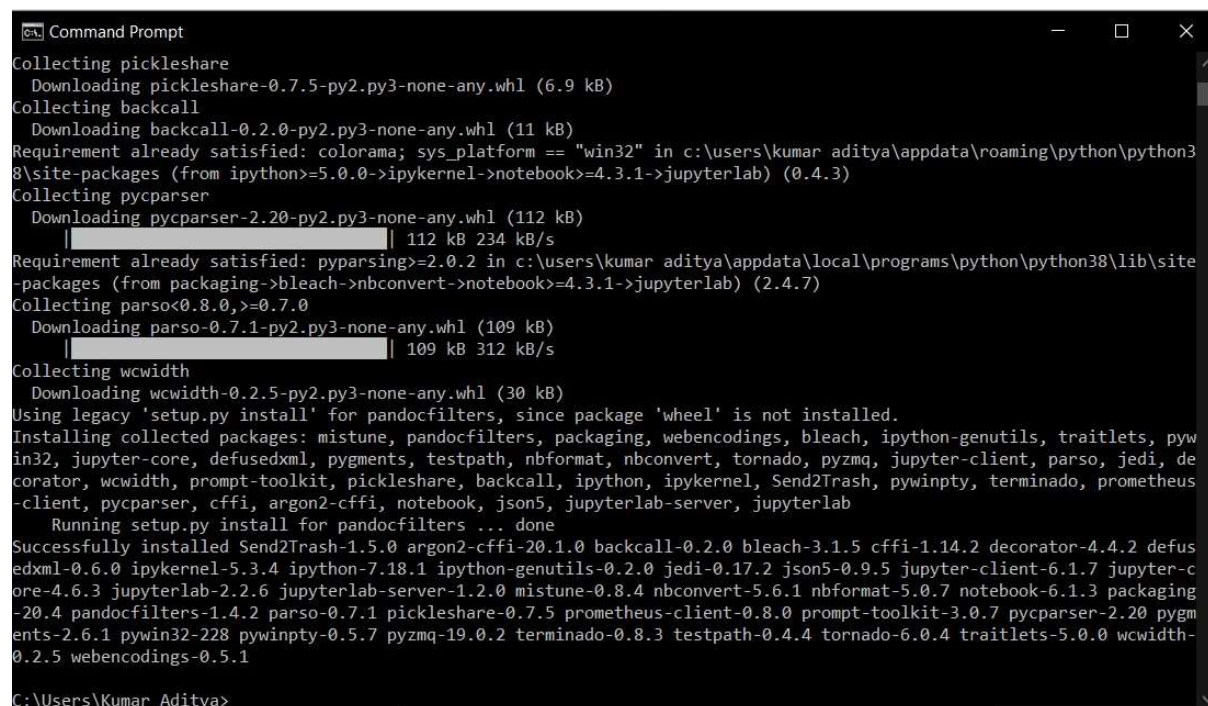
❖ After the completion of setup:



19

### 2.4.1 Opening Jupyter Notebook

You can type following command in cmd to open notebook.

jupyter notebook



After this your notebook will be opened in your default browser:

## 2.5 Google Colaboratory

With Colab you can harness the full power of popular Python libraries to analyse and visualize data. To edit the code, just click the cell and start editing. It is one of the best alternatives for Jupyter Notebook.



## 2.6 INSTALLING PYTHON MODULES:

```
python -m pip install SomePackage
```

# CHAPTER 3

# METHODS AND MATERIAL USED

## 3.1 MATERIAL USED
- Jupyter Notebook/ Google Colaboratory

## 3.2 METHODS
## 3.2.1 DESIGNING THE PROJECT

Confirmed, deaths and recovered are kept in different CSV files. That makes difficult for plotting them in the same data visualization.

This COVID-19 Data processing runs the following steps:
1. Download raw CSV dataset from JHU CSSE public Github page.
2. Load raw CSV dataset and extract the common date list.
3. Merges the raw confirmed, deaths, and recovered CSV data into one DataFrame.
4. Performs data cleanings due to missing values, wrong datatypes and cases from cruise ships.
5. Data Aggregation: Add an active case column *Active*, which is calculated by active_case = confirmed — deaths — recovered. Aggregate data into Country/Region wise and group them by Date and Country/Region. After that, add day wise New cases, New deaths and New recovered by deducting the corresponding cumulative data on the previous day.

There are 3 tasks we would like to do
1. Converting Date from string to datetime
2. Replacing missing value NaN
3. Coronavirus cases reported from 3 cruise ships should be treated differently.

Importing Libraries:

```
[1]  import numpy as np
     import pandas as pd
     import urllib.request
     import altair as alt
     import seaborn as sns
     import matplotlib
```

Getting data from URL:

```
[2]  url1 = 'https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/time_series_covid19_confirmed_global.csv'
     url2 = 'https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/time_series_covid19_deaths_global.csv'
     url3 = 'https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/time_series_covid19_recovered_global.csv'
     urllib.request.urlretrieve(url1,'time_series_covid19_confirmed_global.csv')
     urllib.request.urlretrieve(url2,'time_series_covid19_deaths_global.csv')
     urllib.request.urlretrieve(url3,'time_series_covid19_recovered_global.csv')
     confirmed_df = pd.read_csv('time_series_covid19_confirmed_global.csv')
     deaths_df = pd.read_csv('time_series_covid19_deaths_global.csv')
     recovered_df = pd.read_csv('time_series_covid19_recovered_global.csv')
```

Initial Data:



Confirmed, death and recovered are in separate DataFrame. That makes it difficult to visualize data together.

Merging Data:

```
[12]  full_table = confirmed_df_long.merge(
          right=deaths_df_long,
          how='left',
          on=['Province/State', 'Country/Region', 'Date', 'Lat', 'Long']
      )
      # Merging full_table and recovered_df_long
      full_table = full_table.merge(
          right=recovered_df_long,
          how='left',
          on=['Province/State', 'Country/Region', 'Date', 'Lat', 'Long'])
```

Merged Data:

- Confirmed, Deaths, Recovered and Active are cumulative data.

- New cases, New deaths and New Recovered are day wise data.

- This DataFrame is ordered by Date and Country/Region.

```
[13] full_table.tail()
```

| | Province/State | Country/Region | Lat | Long | Date | Confirmed | Deaths | Recovered |
|---|---|---|---|---|---|---|---|---|
| 59845 | NaN | West Bank and Gaza | 31.952200 | 35.233200 | 9/2/20 | 23875 | 162 | 15483.0 |
| 59846 | NaN | Western Sahara | 24.215500 | -12.885800 | 9/2/20 | 10 | 1 | 8.0 |
| 59847 | NaN | Yemen | 15.552727 | 48.516388 | 9/2/20 | 1976 | 571 | NaN |
| 59848 | NaN | Zambia | -13.133897 | 27.849332 | 9/2/20 | 12415 | 292 | 11494.0 |
| 59849 | NaN | Zimbabwe | -19.015438 | 29.154857 | 9/2/20 | 6638 | 206 | 5250.0 |

Drop columns which are not required:

```
[32] full_table.drop('Province/State', axis=1)
     full_table.drop('Lat', axis=1)
     full_table.drop('Long', axis=1)
```

Plotting graphs:

For simplicity, let's use Python Data Visualization library Altair to create some simple visualizations. Altair is a declarative statistical visualization library for Python, based on Vega and Vega-Lite. Altair offers a powerful and concise visualization grammar that enables you to build a wide range of statistical visualizations quickly.

# CHAPTER 4

## RESULT AND DISCUSSION

## RESULT

Now we have connected everything together and now we will try our data by calling the functions.

➕ Stats for India:





➕ Density Graphs:



Total Death vs Total Recovered

Total Active vs Total Recovered



Total Active vs Total Death

⬇ Comparing different countries (India, Russia, Germany) on Daily New Cases:

**DISCUSSION**

Different government of different countries are working hard to control the spread of COVID -19. While US (Population: 32.82 crores (2019)) saw huge spikes in deaths during the start of this pandemic (approx. 2000 per day) while the country with much more populations like India (135.26 crores (2018)) have only recorded 2000 deaths one day and an average of 1000 deaths per day. Aim of this project is to process and help to analyze effect of measures by different countries.

I have used Python as my basic language because it is easy to understand and very efficient and capable of doing computational work. Libraries like Matplotlib, Seaborn and Alter are very much capable of visualizing data efficiently and libraries like pandas are efficient for manipulating data.

# CHAPTER 5

**WHY DATA PROCESSING?**

The main reason behind data processing is that data almost never comes in a form that is ready for us. A large amount of time spent on a data science project is on manipulating data.

Importance of data processing includes increased productivity and profits, better decisions, more accurate and reliable. Further cost reduction, ease in storage, distributing and report making followed by better analysis and presentation are other advantages. The need to process data is now widely realized and reflected in every field of work. Let the work be done in a business atmosphere or for educational research purpose, data management systems are used by every business. It is a multidimensional process which is involved in almost every field of human life. Generally speaking, the term "Data Processing" is used where you have to collect innumerable data files from different sources. You have to arrange them in a way that can be practically beneficial for the purpose you have gathered all that material. It is a task of synchronizing collected data from different sources and convert it to an organized form. This makes it easy to understand and retrieve the specific information anytime. There are various data processing methods which include manual data processing, mechanical data processing and electronic data processing. Data processing is one of the most important daily tasks especially when dealing with big data and performing data mining. All those fields where we can expect a huge data available to settle down like education, banking or transportation now realizes the importance of data processing.

The invention of computer technology was one of the most important events of all time. With the improvement of computer technology and ease of use, it has become a popular technology in the hands of many. Data processing has also become popular with computer systems making it easier to be handled. In the current times of multiple industries ruling the economy of the many countries of the world, data processing is a field that has numerous applications in most fields like business, education, healthcare, research and more. The importance is increasing with the increase in advancement in areas like data science, machine learning, artificial intelligence, data quality and data security etc.

The invention of computer technology was one of the most important events of all time. With the improvement of computer technology and ease of use, it has become a popular technology in the hands of many. Data processing has also become popular with computer systems making it easier to be handled. In the current times of multiple industries ruling the economy of the many countries of the world, data processing is a field that has numerous applications in most fields like business, education, healthcare, research and more. The importance is increasing with the increase in advancement in areas like data science, machine learning, artificial intelligence, data quality and data security etc.

Data is being collected by almost everyone either knowingly or unknowingly. Collection of data is the first step but processing of data is another vital activity. Companies, institutes & various groups all over the world are engaged in the work of data processing. While talking about the importance of data processing it is equally important to be aware about the related aspects right starting from the methods of data collection, data processing, data processing cycle, information processing cycle, methods of data processing, types of data processing, data presentation and analysis till the data management best practices.


**SPEED, ACCURATE AND RELIABLE**

Data processing is important to make sure that all of the work that is done through that collected facts and figures is done quite speedily and without making any errors. When a data is collected and figured through computers, there are no or negligible chance of errors. It is almost guaranteed that the further processes will be done with maximum possible accuracy. If the input data is accurate then the output is always accurate. Processing can be done at a greater speed and with higher accuracy of right set and combination of software's are used. Another importance of data processing is a major advantage when working in a competitive environment. It is not uncommon to have access to same data. Data and information with higher degree of accuracy is more reliable. Predictive modelling, data cleaning, data validation batch processing is necessary for accurate data.

**COST IS REDUCED**

Data once collected acts as asset for any group and having it stored provides easy access to when required. This eliminates the need to collect data again and again. Moreover, it is very easy and convenient to make copied of the stored data when stored in digital form. Sending or transferring the data is also much easier. This directly helps in cost reduction. The cost or loss which a company might incur because of lack of information is also drastically reduced. This is so as processed data enables it to take a wise and informed decision thus again saving on huge cost.

**Other benefits and merits of Data Processing are:**

1. Data processing makes it easier to validate actions and changes and transactions easily and reduce dependence on computational power for collecting them on demand from a basic form.
2. Insurance claims can be easily handled and settled with properly processed data and make it time-saving for the police authorities as well. Keeping and managing health records, creating electronic health records is now possible due to batch processing, powerful & reliable data warehouses.
3. Data processing can also be made to include image processing and make it easier to present any data to users in a readable format that is liked by them.
4. Invoices can be easily generated for services which have been used and make the customer experience better.
5. Data processing in the form of word processing can help in making documents which are readable and likeable by readers and be made even more engaging.

# CHAPTER 6

## CONCLUSION & FUTURE WORK

### SUMMARY:

The intension of this project is to process different data sets of COVID – 19 such as confirmed, recovered and death and plots statics graph to help analyze the approach and its effect on different countries.

### FUTURE WORK:

In future projects I will improve the design of my program by giving User Interface and more graph options. Also, would put more extra features like prediction modelling of countries.

### CONCLUSION:

Open the Jupyter Notebook and run all the cell. When prompted enter country names or a set of country name as necessary for the analysis.

# CHAPTER 7

## SOURCE CODE:

```python
[1] import numpy as np
    import pandas as pd
    import urllib.request
    import altair as alt
    import seaborn as sns
    import matplotlib
```

```python
[2] url1 = 'https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/time_series_covid19_confirmed_global.csv'
    url2 ='https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/time_series_covid19_deaths_global.csv'
    url3= 'https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/time_series_covid19_recovered_global.csv'
    urllib.request.urlretrieve(url1,'time_series_covid19_confirmed_global.csv')
    urllib.request.urlretrieve(url2,'time_series_covid19_deaths_global.csv')
    urllib.request.urlretrieve(url3,'time_series_covid19_recovered_global.csv')
    confirmed_df = pd.read_csv('time_series_covid19_confirmed_global.csv')
    deaths_df = pd.read_csv('time_series_covid19_deaths_global.csv')
    recovered_df = pd.read_csv('time_series_covid19_recovered_global.csv')
```

```python
[3] dates = confirmed_df.columns[4:]
    confirmed_df_long = confirmed_df.melt(
        id_vars=['Province/State', 'Country/Region', 'Lat', 'Long'],
        value_vars=dates,
        var_name='Date',
        value_name='Confirmed'
    )
    deaths_df_long = deaths_df.melt(
        id_vars=['Province/State', 'Country/Region', 'Lat', 'Long'],
        value_vars=dates,
        var_name='Date',
        value_name='Deaths'
    )
    recovered_df_long = recovered_df.melt(
        id_vars=['Province/State', 'Country/Region', 'Lat', 'Long'],
        value_vars=dates,
        var_name='Date',
        value_name='Recovered'
    )
```

```python
[4] full_table = confirmed_df_long.merge(
        right=deaths_df_long,
        how='left',
        on=['Province/State', 'Country/Region', 'Date', 'Lat', 'Long']
    )
    full_table = full_table.merge(
        right=recovered_df_long,
        how='left',
        on=['Province/State', 'Country/Region', 'Date', 'Lat', 'Long'])
```

```python
[5] full_table['Date'] = pd.to_datetime(full_table['Date'])
```

```python
[6] full_table['Recovered'] = full_table['Recovered'].fillna(0)
```

```python
[7] full_table.drop('Province/State', axis=1,inplace=True)
    full_table.drop('Lat', axis=1,inplace=True)
    full_table.drop('Long', axis=1,inplace=True)
```

```python
[8] full_table['Active'] = full_table['Confirmed'] - full_table['Deaths'] - full_table['Recovered']
```

```python
[9] full_grouped = full_table.groupby(['Date', 'Country/Region'])['Confirmed', 'Deaths', 'Recovered', 'Active'].sum().reset_index()
```
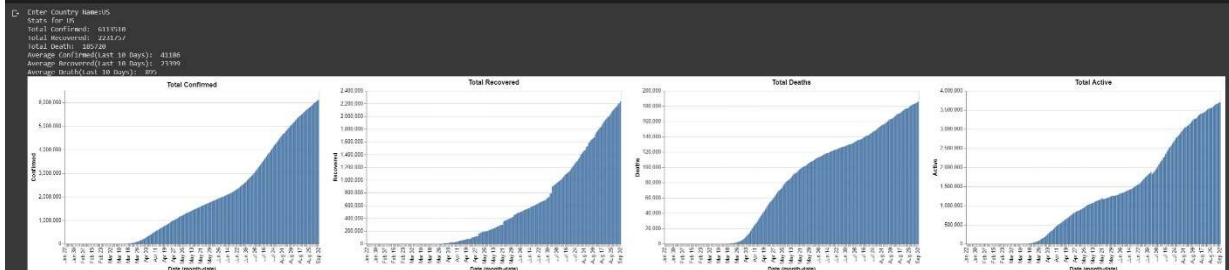
```python
[10] temp = full_grouped.groupby(['Country/Region', 'Date', ])['Confirmed', 'Deaths', 'Recovered']
     temp = temp.sum().diff().reset_index()
     mask = temp['Country/Region'] != temp['Country/Region'].shift(1)
     temp.loc[mask, 'Confirmed'] = np.nan
     temp.loc[mask, 'Deaths'] = np.nan
     temp.loc[mask, 'Recovered'] = np.nan
     temp.columns = ['Country/Region', 'Date', 'New cases', 'New deaths', 'New recovered']
     full_grouped = pd.merge(full_grouped, temp, on=['Country/Region', 'Date'])
     full_grouped = full_grouped.fillna(0)
     cols = ['New cases', 'New deaths', 'New recovered']
     full_grouped[cols] = full_grouped[cols].astype('int')
     full_grouped['New cases'] = full_grouped['New cases'].apply(lambda x: 0 if x<0 else x)
```

```
[11] full_grouped.to_csv('COVID-19-time-series-clean-complete.csv')
     full_grouped = pd.read_csv('COVID-19-time-series-clean-complete.csv', parse_dates=['Date'])
```
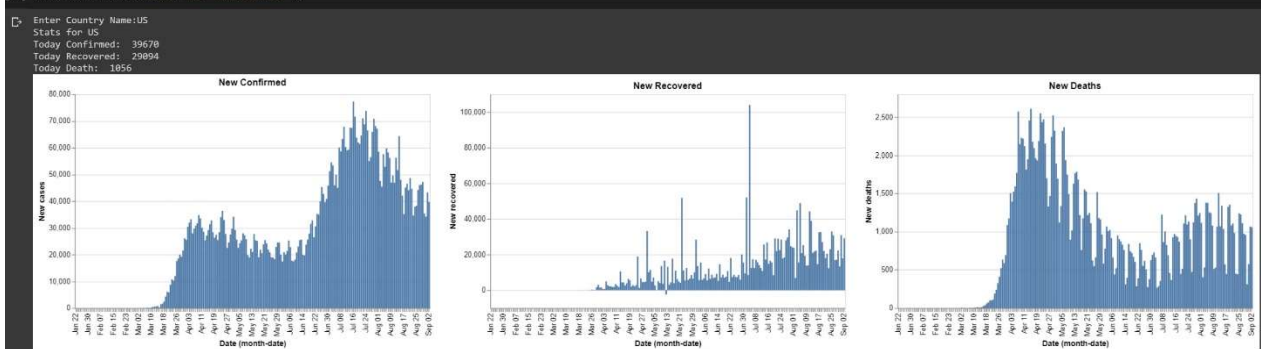
```
[12] def total_country_stats(country_name):
         stat_country = full_grouped[full_grouped['Country/Region'] == country_name]
         base = alt.Chart(stat_country).mark_bar().encode(x='monthdate(Date):O',).properties(width=500)
         red=alt.value('#f54242')
         print("Stats for "+country_name, end="\n")
         data_current=stat_country.tail(1)
         data_10=stat_country.tail(10)
         print("Total Confirmed: ",int(data_current['Confirmed'].mean()))
         print("Total Recovered: ",int(data_current['Recovered'].mean()))
         print("Total Death: ",int(data_current['Deaths'].mean()))
         print("Average Confirmed(Last 10 Days): ",int(data_10['New cases'].mean()))
         print("Average Recovered(Last 10 Days): ",int(data_10['New recovered'].mean()))
         print("Average Death(Last 10 Days): ",int(data_10['New deaths'].mean()))
         return base.encode(y='Confirmed').properties(title='Total Confirmed') | base.encode(y='Recovered').properties(title='Total Recovered') |
         base.encode(y='Deaths').properties(title='Total Deaths') | base.encode(y='Active').properties(title='Total Active')
```

```
[13] def daily_country_stats(country_name):
         stat_country = full_grouped[full_grouped['Country/Region'] == country_name]
         base = alt.Chart(stat_country).mark_bar().encode(x='monthdate(Date):O',).properties(width=500)
         red=alt.value('#f54242')
         print("Stats for "+country_name)
         data_current=stat_country.tail(1)
         print("Today Confirmed: ",int(data_current['New cases'].mean()))
         print("Today Recovered: ",int(data_current['New recovered'].mean()))
         print("Today Death: ",int(data_current['New deaths'].mean()))
         return base.encode(y='New cases').properties(title='New Confirmed') | base.encode(y='New recovered').properties(title='New Recovered') |
         base.encode(y='New deaths').properties(title='New Deaths')
```
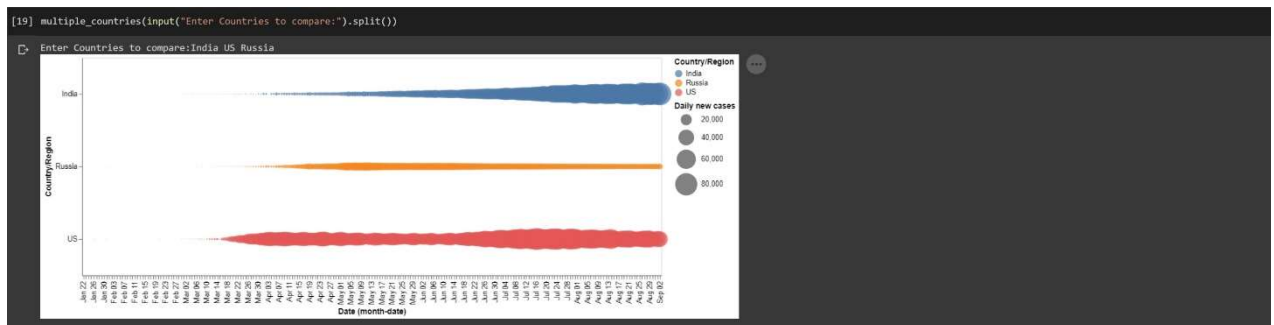




```
[16] def country_jointplots(country_name):
         stat_country = full_grouped[full_grouped['Country/Region'] == country_name]
         df = pd.DataFrame({'Total Death': stat_country['Deaths'], 'Total Recovered': stat_country['Recovered']})
         sns.jointplot('Total Death', 'Total Recovered', data=df,kind = 'kde')
         df = pd.DataFrame({'Total Active': stat_country['Active'], 'Total Recovered': stat_country['Recovered']})
         sns.jointplot('Total Active', 'Total Recovered', data=df,kind = 'kde')
         df = pd.DataFrame({'Total Active': stat_country['Active'], 'Total Death': stat_country['Deaths']})
         sns.jointplot('Total Active', 'Total Death', data=df,kind = 'kde')
```

```
country_jointplots(input("Enter Country Name:"))
```

Enter Country Name:US



```
[18] def multiple_countries(country_list):
        selected_countries = full_grouped[full_grouped['Country/Region'].isin(country_list)]
        return alt.Chart(selected_countries).mark_circle().encode(x='monthdate(Date):O',y='Country/Region',color='Country/Region',size=alt.Size('New cases:Q',scale=alt.Scale(range=[0, 1000]),legend=alt.Legend(title='Daily new c
```

```
[19] multiple_countries(input("Enter Countries to compare:").split())
```

Enter Countries to compare:India US Russia

# CHAPTER 8

**REFERENCES:**

I.      IBM Data Science Professional Coursera

II.      Geeks for Geeks

III.      Stack Over Flow

IV.      Krish Naik Youtube