Statistics For Data Science

**Sampriti Chatterjee (Great Learning)**

# Agenda


greatlearning
*Learning for Life*

1. Why do we need data science?

2. What is Data science?

3. Life cycle of Data science

4. Important statistics terms in data science

5. Install python

6. Python Library: Numpy and Pandas

7. Data manipulation using Numpy and Pandas

8. Data visualization with seaborn and Matplotlib

9. What is machine Learning?

10. Supervised Learning: Logistic Regression

11. Diabetes prediction using Python

# Why do we need Data Science?

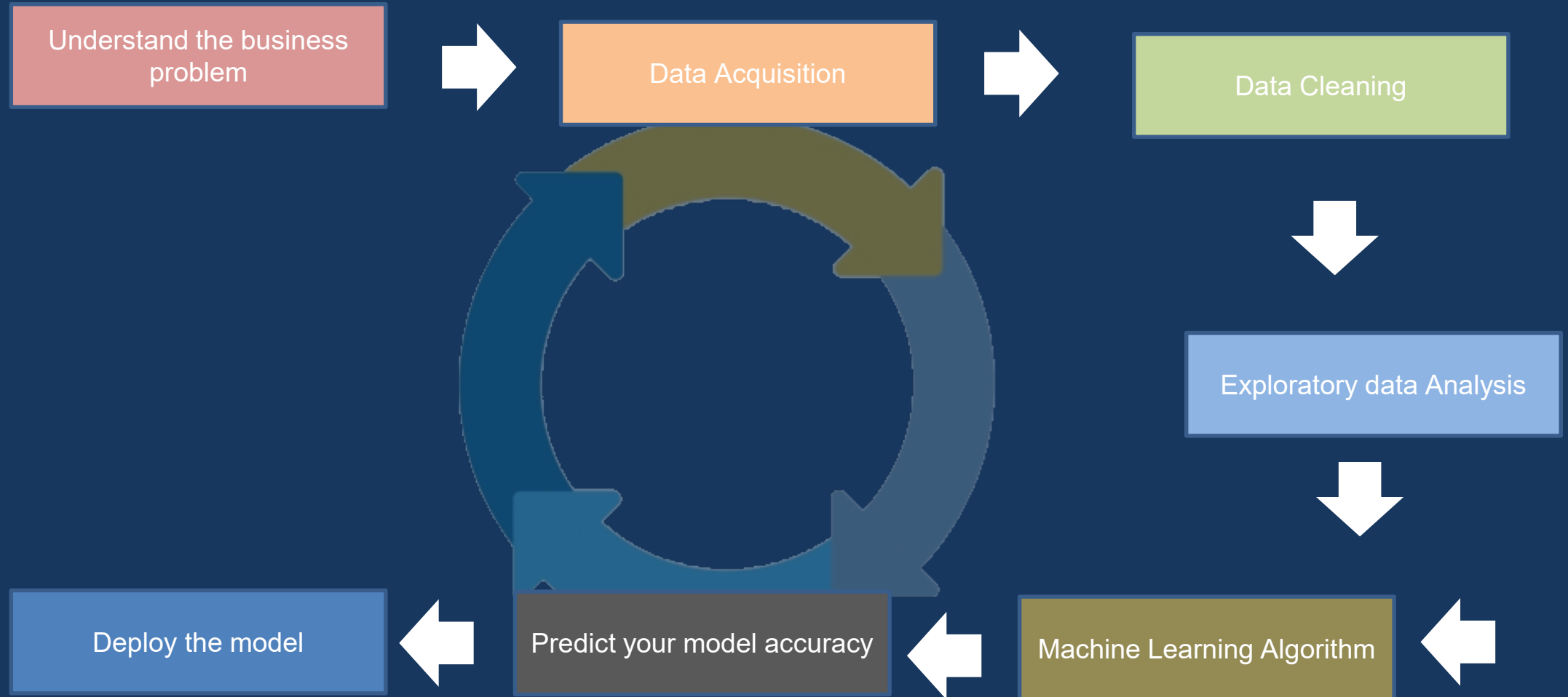How data science is effecting our everyday life?

- In the past, we used to have data in a structured format but now as the volume of the data is increasing, so the number of structured data becomes very less, so to handle the massive amount of data we need data science techniques

- Those data can be used to get the proper business insights and the hidden trends from them.

- These insights helps the organization to predict the Future

- Using data science decision making can be faster and effective

- Helps to reduce the production cost

- Build model based on the data to give the ability to the machine to predicts on its own
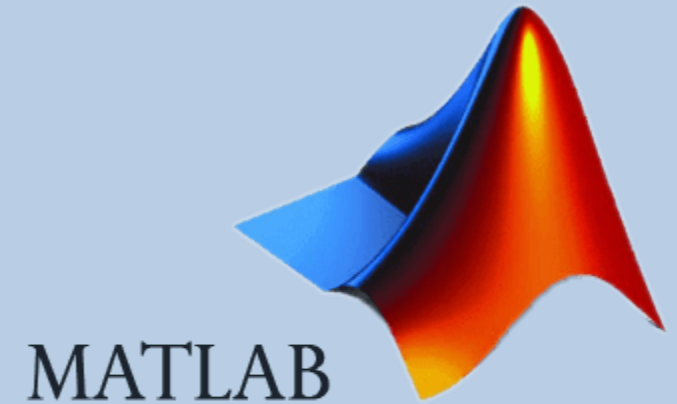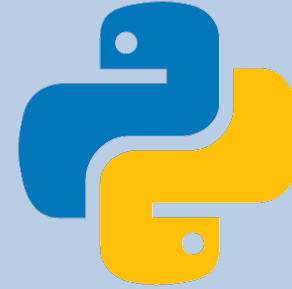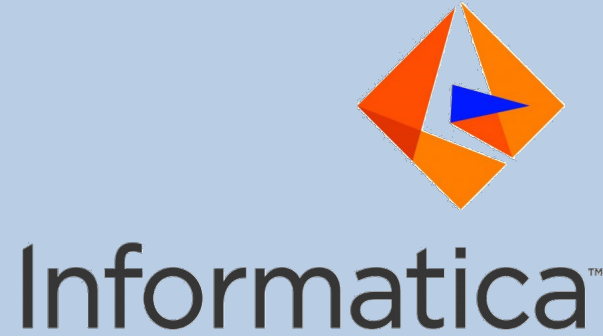
greatlearning
Learning for Life

Data science is a process to get some meaningful information from the massive amount of data. In simple terms, read and study the data to get proper intuitive insights. Data Science is a mixture of various tools, algorithms, and machine learning and deep learning concepts to discover hidden patterns from the raw and unstuctured data

Most Popular Programming Languages For Data Science?

# Important statistics terms in data science

1 What is Statistics?

2 What is population?

3 What is parameter?

4 What is sample?

5 What is mean?

6 Types of analysis in statistics

7 What is Outlier?

8 What is Interquartile Range IQR?

9 What is upper and lower limits in interquartile range

10 What is null hypothesis?

11 What is p value?

**greatlearning**
*Learning for Life*

# What is Statistics?

**Statistics** is a part of integrated applied mathematics which deals with data

**1** It helps to collect data and analyze them properly

**2** With the help of statistics we can read the data and organize them in order to get the hidden information from them

**3** In data science domain statistics concepts are used to process the complex data to get the insights from them using mathematical computations
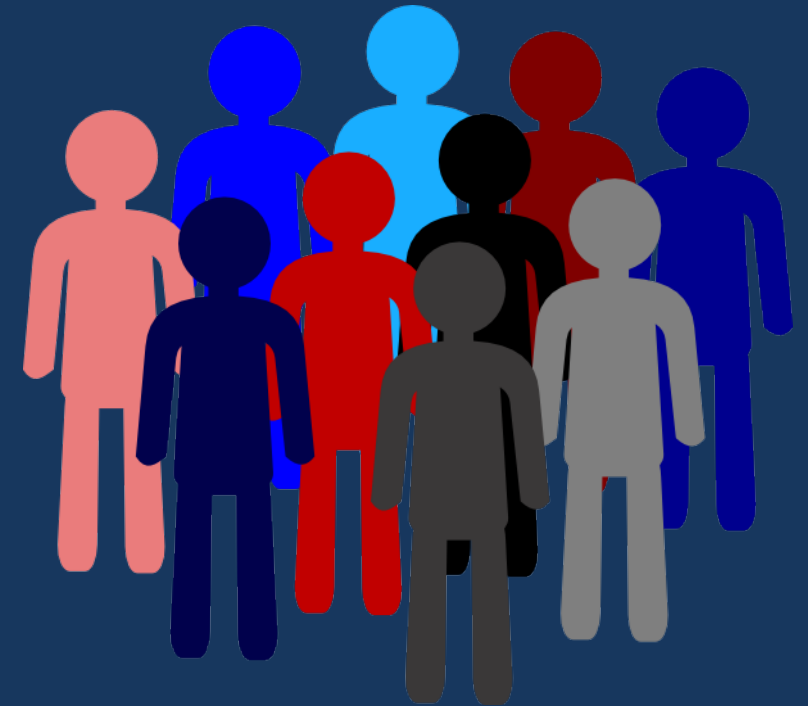
greatlearning
*Learning for Life*

# What is Population?

Population terms in statistics use to refer the total set of observations

**Example:**

Suppose,
If we want to study a diabetes dataset to understand the symptoms and the other factors then the whole dataset is referred as population

# What is Parameter?

Parameters are referred to characteristics which describes the population
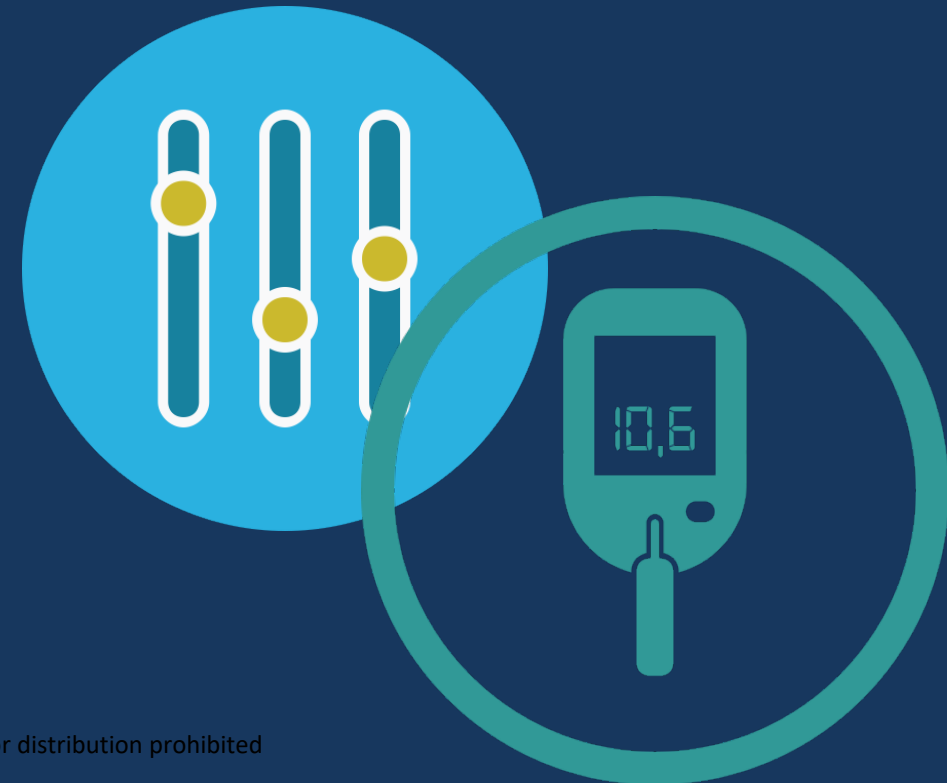
1 Parameters are like average or percentage which helps to describe the entire population

2 Mean and the standard deviation are two common parameters of population

3 Example: Average age for being diabetic is the parameter for whole diabetes dataset population

# What is Sample?

Sample is basically a small part or portion of the large population

**Example:**

Suppose,
From the whole diabetes dataset you picked
100 rows of information to do the analysis, that
100 rows of information will be referred as
**Sample**

Target Population

Sample

greatlearning
*Learning for Life*

**What is Mean?** | Mean term referred as average value of the whole population

**What is Median?** | Median is the middle value of the data when your data is sorted in manner

**What is Mode?** | Mode stands for the most occurring element in the dataset

great**learning**
*Learning for Life*

# Types Of Analysis In Statistics

**Descriptive statistics**

**Inferential Statistics**

It helps to describe the data in mathematical or graphical way

Inferential statistics split the data into samples and applies probability to arrive to theconclusion

greatlearning
*Learning for Life*

# What is Outlier?

Outliers in the dataset are referred as unusual value which can distort and violate statistical analysis

**1**    Outliers are basically experimental errors in the data

**2**    Some outliers are good for the dataset to detect anomaly like: detecting fraud transaction

**3**    It effects the mean and the standard deviation of the data and most of the machine learning technique does not perform good with outliers

# What is Interquartile Range IQR?

Interquartile range divides the dataset into quartiles to measure the variability and the spread of the dataset

**1** Splits the data into 4 equal part in sorted manner

**2** Q1, Q2, Q3 are called first, second and third quartiles:

- Q1 → 25th percentile of the dataset

- Q2 → 50th percentile of the dataset

- Q3 → 75th percentile of the dataset

**Formula: IQR → Q3 – Q1**

greatlearning
*Learning for Life*

# What is upper and lower limits in interquartile range

Lower and upper limit in the interquartile basically the range where data points lie

**1** **Formula to find the lower limit:**
**Lower_limit = Q1 - 1.5 IQR**

**2** **Formula to find the upper limit:**
**Upper_limit = Q3 + 1.5 * IQR**

greatlearning
Learning for Life

# What is p value?

p value is used to support or reject the null hypothesis or the assumption

1. P value is basically the strong evidence to reject the null hypothesis

2. If p value is less than 0.05 then we accept the null hypothesis

**Python is a popular high level, object oriented and interpreted language**



High level

Interpreted

Object oriented

This is the site to install Python -> https://www.python.org/downloads/

# Popular IDE for Python: Pycharm

Site to install Python ->
https://www.jetbrains.com/pycharm/download/#section=mac

Anaconda installation site->
https://www.anaconda.com/products/individual

# Popular IDE for Python: Google colab

Google collaboratory link->
https://colab.research.google.com/notebooks/intro.ipynb

Getting started with Python

# Python Libraries

Data manipulation is a technique which allows to transform, extract, and filter your data efficiently with less time.

Main two python libraries are used to manipulate the data

NumPy

Pandas

![greatlearning - Learning for Life]

**Numpy stands for Numerical Python and it is used to perform mathematical and logical operations on arrays**

**1**    Numpy is a python library

**2**    Install Numpy:  !pip install numpy

**3**    Import the Library: import numpy as np

# Data manipulation using Pandas

**Pandas is a popular data manipulation and analysis library in python which is based on Numpy**

1 — Python is a python library built on top of Numpy

2 — Install Numpy:  !pip install pandas

3 — Import the Library: import pandas as pd

Pandas

Demo on Numpy and pandas

# greatlearning
*Learning for Life*

Panda's data frame is a two-dimensional data structure which is aligned in a tabular fashion with rows and columns

What is a dataframe?

```
In [9]: import pandas as pd

        pd.DataFrame({"Name":['Bob','Sam','Anne'],"Marks":[76,25,92]})

Out[9]:
              Name    Marks
        0     Bob        76
        1     Sam        25
        2     Anne       92
```

head()

shape()

describe()

tail()

# Dropping Columns

**greatlearning**
*Learning for Life*

```
iris.drop('Sepal.Length',axis=1)
```

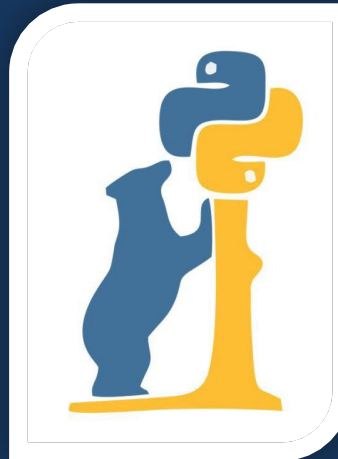|   | Sepal.Length | Sepal.Width | Petal.Length | Petal.Width | Species |
|---|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | setosa |

→

|   | Sepal.Width | Petal.Length | Petal.Width | Species |
|---|---|---|---|---|
| 0 | 3.5 | 1.4 | 0.2 | setosa |
| 1 | 3.0 | 1.4 | 0.2 | setosa |
| 2 | 3.2 | 1.3 | 0.2 | setosa |
| 3 | 3.1 | 1.5 | 0.2 | setosa |
| 4 | 3.6 | 1.4 | 0.2 | setosa |

# Dropping Rows

```
iris.drop([1,2,3],axis=0)
```

| | Sepal.Length | Sepal.Width | Petal.Length | Petal.Width | Species |
|---|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | setosa |

→

| | Sepal.Length | Sepal.Width | Petal.Length | Petal.Width | Species |
|---|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | setosa |
| 5 | 5.4 | 3.9 | 1.7 | 0.4 | setosa |
| 6 | 4.6 | 3.4 | 1.4 | 0.3 | setosa |
| 7 | 5.0 | 3.4 | 1.5 | 0.2 | setosa |

Machine Learning to build the Model

# Traditional Vs Machine Learning

**Traditional Programming**
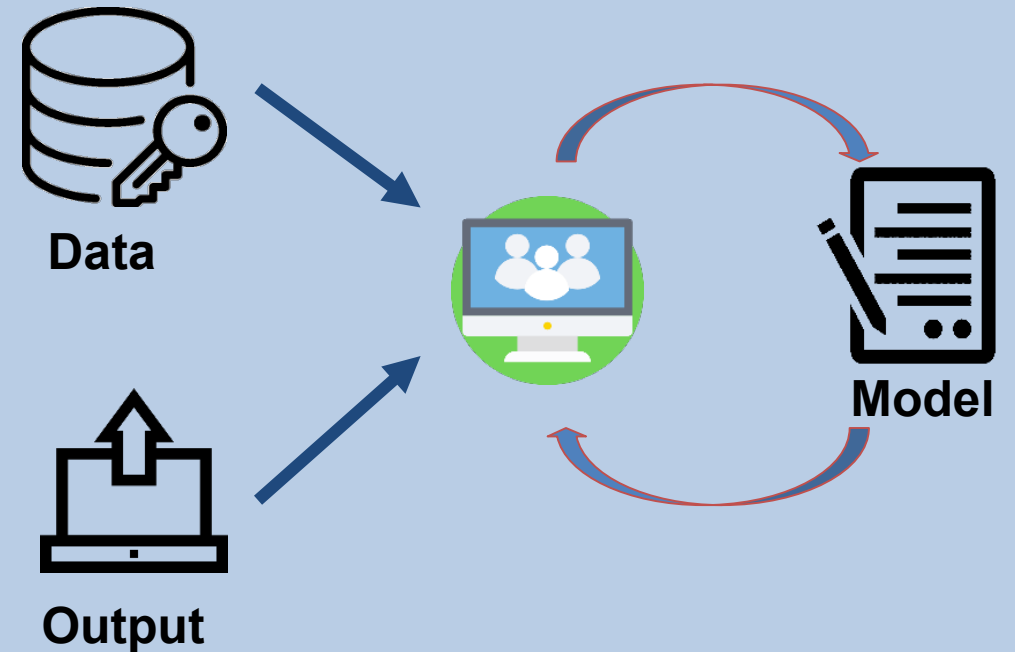
**Machine Learning**

Data

Program

Output

Data

Output

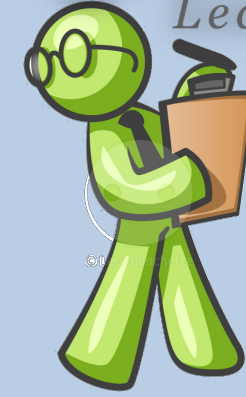Model

# Types Of Machine Learning

greatlearning
Learning for Life

Supervised Learning

Unsupervised Learning

Reinforcement Learning

Environment

State

Action

Reward

Agent

Supervised learning works as a supervisor or teacher. Basically, In supervised learning, we teach or train the machine with labeled data (that means data is already tagged with some predefined class). Then we test our model with some unknown new set of data and predict the level for them

- Learning from the labelled data and applying the knowledge to predict the

  label of the new data(test data), is known as *Supervised Learning*

- *Types of Supervised Learning:*

  - **Linear Regression**

  - **Logistic regression**

  - **Decision Tree**

  - **Random Forest**

  - **Naïve Bayes Classifier**

**greatlearning**
*Learning for Life*
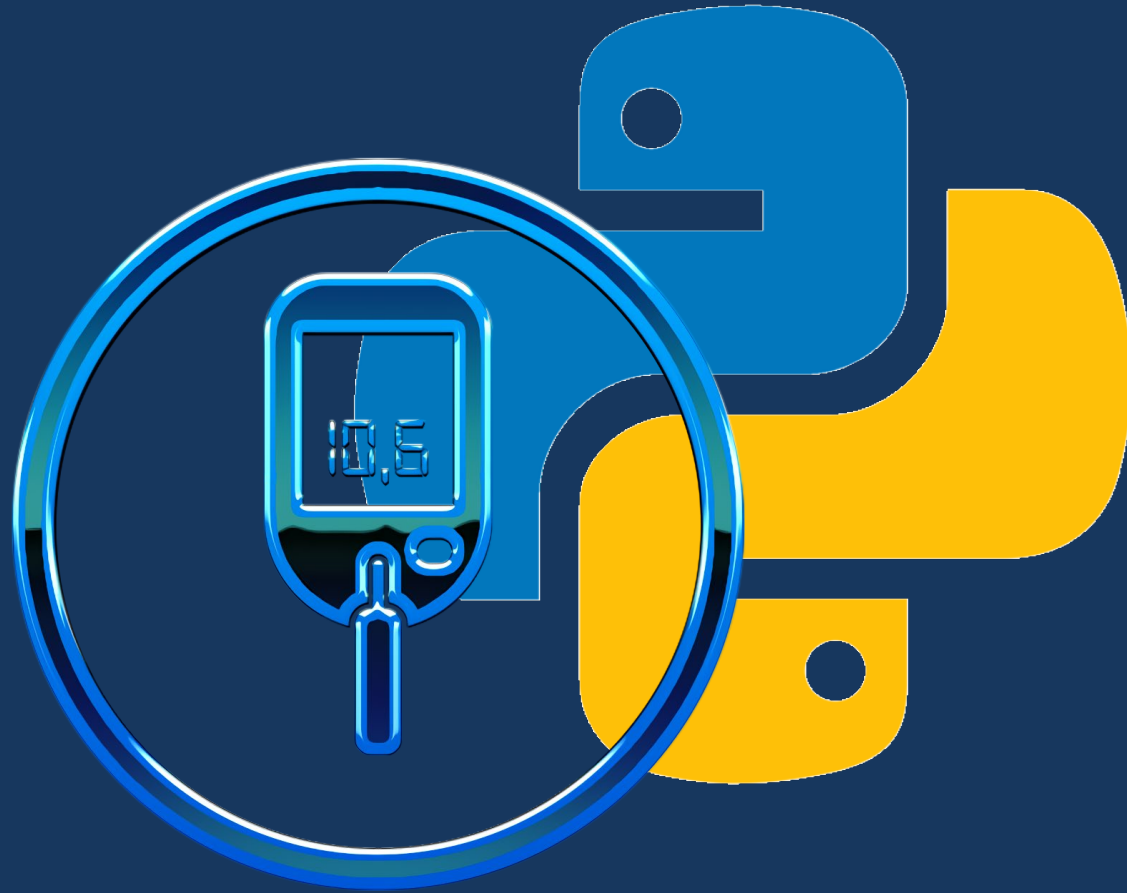
Logistic regression is also a part of supervised learning classification algorithm. It is used to predict the probability of a target variable and the nature of target or dependent variable is discrete, so for the output there will be only two class will be present

- The dependent variable is binary in nature so that can be either 1 (stands for success/yes) or 0 (stands for failure/no).

- Logistic regression is also known as sigmoid function

- *Sigmoid function = 1 / (1 + e^-value)*

# Diabetes Prediction using Python

# Thank You