# EDA for the Iris Flower Dataset



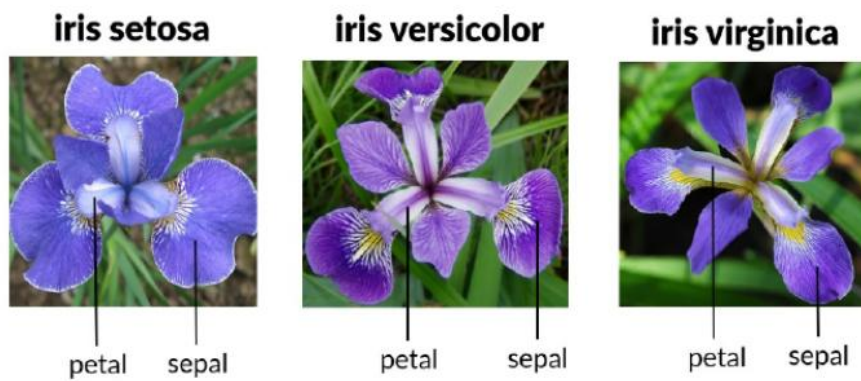iris setosa     iris versicolor     iris virginica

petal   sepal     petal   sepal     petal   sepal

**ANKIT RAJ**
**16/09/2023**

# INTRODUCTION

Every machine learning project begins by understanding what the data and drawing the objectives. While applying machine learning algorithms to your data set, you are understanding, building and analyzing the data as to get the end result.
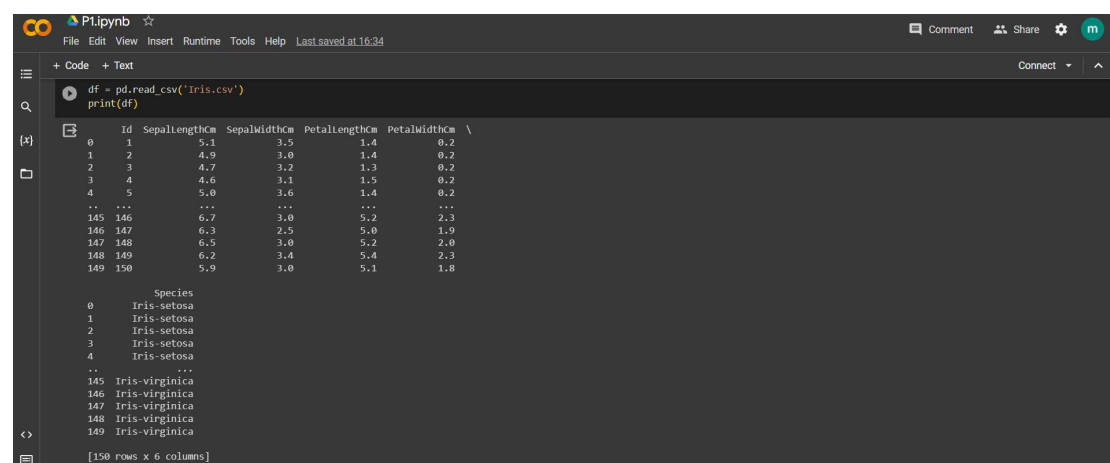
Following are the steps involved in creating a well-defined ML project:
1] Understand and define the problem
2] Prepare the data
3] Explore and Analyse the data

1.1 Problem statement This data set consists of the physical parameters of three species of flower - Versicolor, Setosa and Virginica. The numeric parameters which the dataset contains are Sepal width, Sepal length, Petal width and Petal length. In this data we will be predicting the classes of the flowers based on these parameters. The data consists of continuous numeric values which describe the dimensions of the respective features.We will be training the model based on these features.

1.2 Data Prepare

It has been created Ronald Fisher in 1936. It contains the petal length, petal width,sepal length and sepal width of 150 iris flowers from 3 different species. Variables present in given dataset are SepalLengthCm, SepalWidthCm, PetalLengthCm, PetalWidthCm, Species.
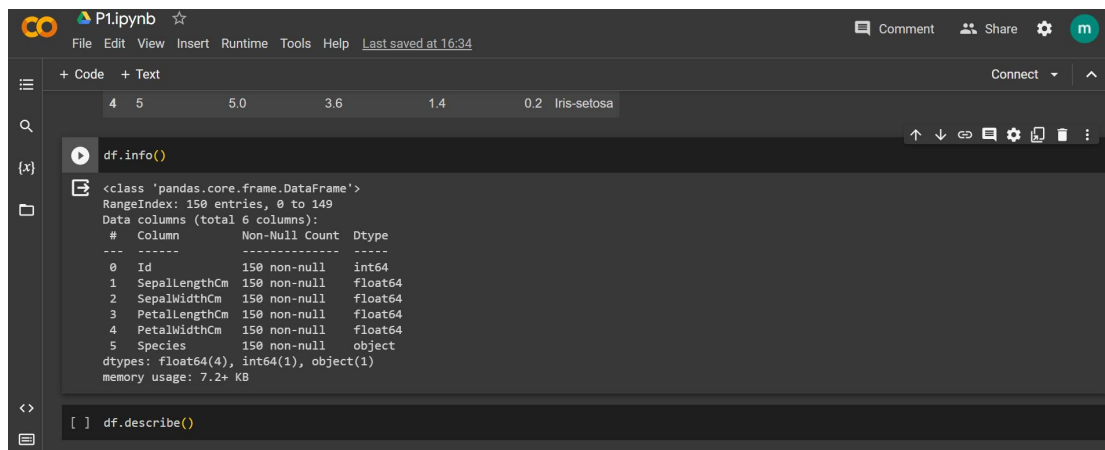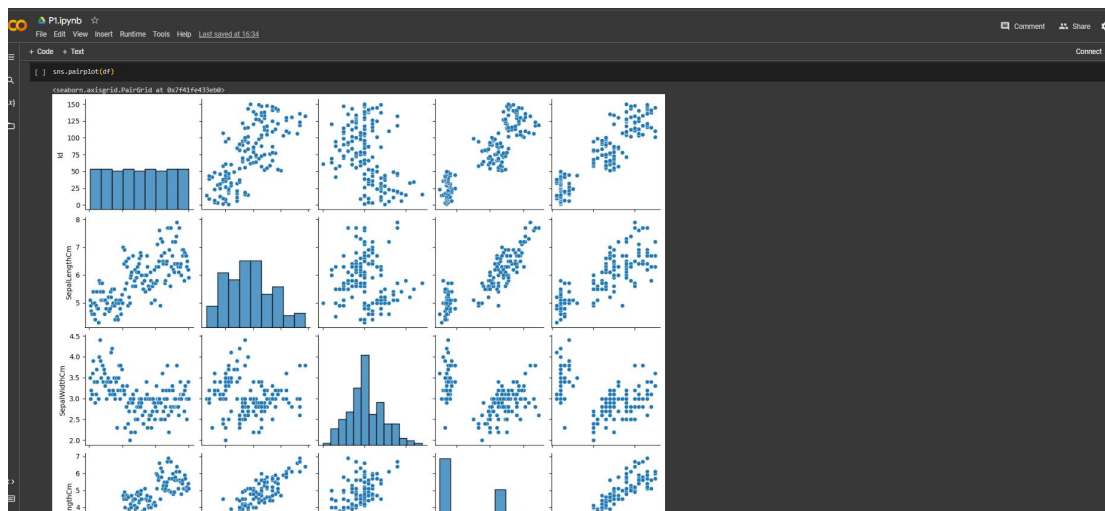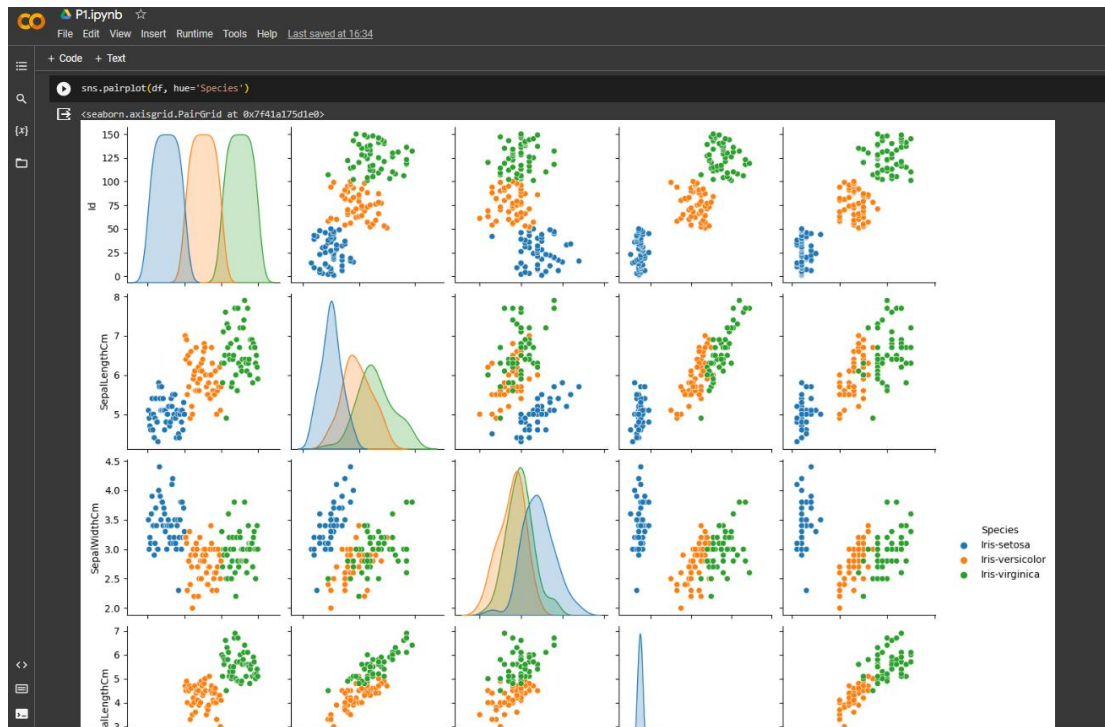
Now, View the info of the data frame that contains details like the count of non-null variables and the column's datatype along with the column names. It will also show the memory usage.

# 2. Methodology

2.1 Pre Processing Any predictive modeling requires that we look at the data before we start modeling. However, in data mining terms looking at data refers to so much more than just looking. Looking at data refers to exploring the data, cleaning the data as well as visualizing the data through graphs and plots. This is often called as Exploratory Data Analysis.

2.1.1 Exploratory Data Analysis In exploring the data we have,If there are any missing values, then modify them before using the dataset. Formodifying you can use the fillna() method. It will fill null values.

+ Code   + Text

```
sns.countplot(x='Species', data=df)
```

<Axes: xlabel='Species', ylabel='count'>



```
sns.histplot(data=df, x='SepalLengthCm')
```

<Axes: xlabel='SepalLengthCm', ylabel='Count'>

```
import statistics
```

```
df.mean()
```

<ipython-input-56-c61f0c8f89b5>:1: FutureWarning: The default value of numeric_only in DataFrame.mean is deprecated. In a future version, it will default to False. In addition, specifying 'numeric_only=None' is deprecated. Select only valid columns o
  df.mean()
```
Id              75.500000
SepalLengthCm    5.843333
SepalWidthCm     3.054000
PetalLengthCm    3.758667
PetalWidthCm     1.198667
dtype: float64
```

```
df.median()
```

<ipython-input-57-6d467abf240d>:1: FutureWarning: The default value of numeric_only in DataFrame.median is deprecated. In a future version, it will default to False. In addition, specifying 'numeric_only=None' is deprecated. Select only valid columns o
  df.median()
```
Id              75.50
SepalLengthCm    5.80
SepalWidthCm     3.00
PetalLengthCm    4.35
PetalWidthCm     1.30
dtype: float64
```

```
df.std()
```

<ipython-input-62-ce97bb7eaef8>:1: FutureWarning: The default value of numeric_only in DataFrame.std is deprecated. In a future version, it will default to False. In addition, specifying 'numeric_only=None' is deprecated. Select only valid columns o
  df.std()
```
Id              43.445368
SepalLengthCm    0.828066
SepalWidthCm     0.433594
PetalLengthCm    1.764420
PetalWidthCm     0.763161
dtype: float64
```

## Conclusion

Flower classification is a very important, simple, and basic project for any machine learning student. Every machine learning student should be thorough with the iris flowers dataset.