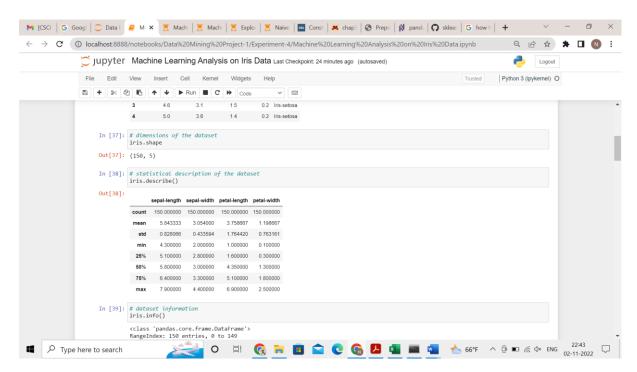# Project-1

**Experiment-4:**

**Iris data:**

**Context:** The iris dataset is about finding the presence of diabetes in a person depending upon the attributes such as {sepal length, sepal width, petal length, petal width}.

**Content:** The input dataset has 150 records with 5 columns where first 4 columns are the features of the dataset and the last column has the label. The output label is a categorical datatype and the dataset has outcome with 50 Iris-setosa, 50 Iris-versicolor and 50 Iris-virginica. All features are of continuous data type.

**Missing values or null data points:** There is no missing or null data in the given dataset.

**Feature:** From the histogram data grouped between individual feature and outcome, all features have variations in their plots compared to other features, hence, all these features have significant impact while training the model.

**Feature Description:**



**Testing Strategy:** Stratified K Fold with 10 splits is the testing strategy deployed.

**Classification accuracy of algorithms:**

Using Train/ Test split:

| Model | Accuracy |
|-------|----------|
| KNN | 1.000000 |
| DB | 0.973684 |
| GNB | 0.315789 |
| BNB | 0.973684 |

Using K-Fold cross validation:

| Model | Accuracy |
|-------|----------|

# Project-1

| KNN | 0.966667 |
|-----|----------|
| DB  | 0.953333 |
| GNB | 0.333333 |
| BNB | 0.953333 |

KNN has the highest accuracy in both the train/ test split and K-Fold validation.
Hence, KNN would be a better model in the case of Iris data.