

Assignment 5: Naive Bayes Model

(Total Marks: 15)

Problem Statement:

A floriculture research team X is studying the use of multiple measurements to distinguish three different iris flower species. The dataset contains a set of 150 records under five attributes: sepal length, sepal width, petal length, petal width and species (see Fig. 1). Develop a Naive Bayes classifier that classifies the species according to the above measurements.

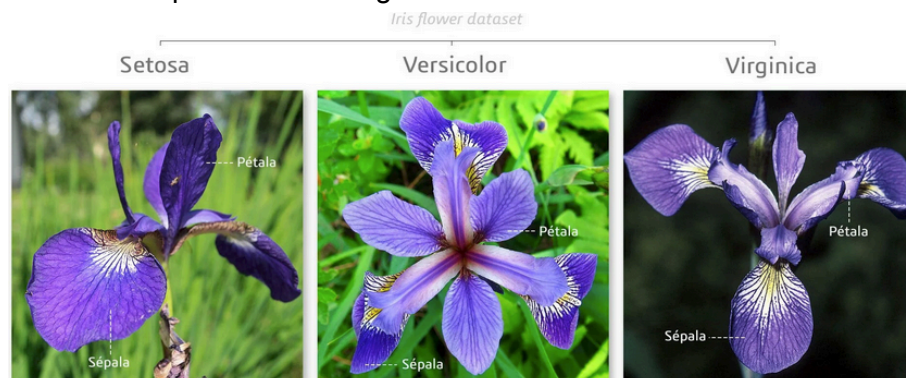


Figure 1: Different iris flower species and their attributes

Implementation: [5 Marks]

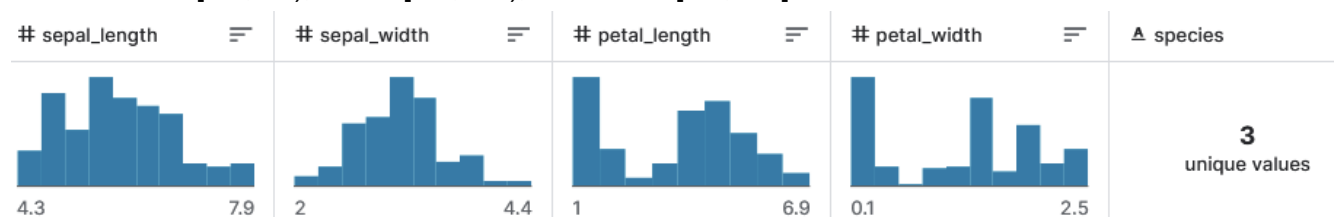
- Implementation of Naive Bayes classifier (NB-CLS) from scratch (without using builtin functions). Evaluate the model using Percentage Accuracy.

****Implement [NB-CLS] from scratch.** You may make use of the numpy library to perform basic operations (e.g., sorting).

****In general,** you may use libraries to process and handle data.

****DO NOT** perform feature scaling before feeding the data in your model.

**** As each attribute spans numerous values (see below), divide each attribute value into K equally wide bins spanning the lower and higher values. For example, sepal_length bins for K =3 are as follows: Bin 1: [4.3,5.5); Bin 2: [5.5, 6.7); and Bin 3: [6.7, 7.9].**



**** To save computation,** you may pre-compute required conditional probabilities and store that in a matrix.

Experiments: [5+3=8 Marks]

The dataset will be split into Train:Test with 80:20 ratio. Pl shuffle the data before splitting.

- Experiment 1:** Report the effect of varying the number of bins K in [NB_CLS] on Test data. Choose K values from [2, 3, 5]. Plot Percentage Accuracy vs K. Find the best value of the hyperparameter K.
- Experiment 2:** Add noise to only a fraction of the training data: consider separately 10%, 40%, 80%, 90% of the training data for noise addition. Choose a normal distribution with zero mean and standard deviation 2.0. Next, design a Naive Bayes using the optimal K found in the earlier experiment. How does the performance vary as compared to that of the noiseless case (Experiment 1)?

(For noise, one may use: `numpy.random.normal(loc=mean, scale=std_dev, size=train_data.shape)` with seed)

Report your observations with appropriate explanations.

Datasets:

This dataset comprises three iris species with 50 samples each as well as some properties about each flower. You can find the dataset [here](#).

- ID: Identification number of the flower
- Sepal length: Length of sepal in cm (in real numbers)
- Sepal Width: Width of flower sepal in cm (in real numbers)
- Petal length: Length of flower petal in cm (in real numbers)
- Petal Width: Width of flower petal in cm (in real numbers)
- Species: Three iris flower species (iris-setosa, iris-versicolor, and iris-virginica)

Problem: Predict the species of an iris flower

Submission:

A .zip file containing the python source code and a PDF report file. The final name should follow the template: <Assign-No>_<Your Roll No>.zip. For example, if your roll no is 15CE30021, the filename for Labassign-4 will be: [LabAsgn-4_15ce30021.zip](#)

1. A **single python code (.py)** containing the implementations of the models and experiments with comments at function level. The first two lines should **contain your name and roll no**.
2. A report [PDF] containing **[2 Marks]**
 - a. Experiment 1: Plot of Percentage Accuracy vs K. Also mention the best choice for the K and the corresponding percentage accuracy.
 - b. Experiment 2: Report the performance at different noise levels. Comment on the robustness of NB-CLS to noise in the training dataset.

Responsible TAs: Please write to the following TAs for any doubt or clarification regarding Assignment 4.

Soumyadipto Banerjee - soumyadiptobnrj071@gmail.com

Deadline: The deadline for submission is **12th February (Monday), 11:55 PM, IST**. Irrespective of the time in your device, once submission in moodle is closed, no request for submission post-deadline will be entertained. No email submission will be considered. So, it is suggested that you start submitting the solution at least one hour before the deadline.

Plagiarism policy: Binary marking (two parties)