

# CSCM35, CSLM35 Big Data and Data Mining

by Dr. Jingjing Deng

Released on 27th Feb 2020

The assignment consists of one task that are designed to be completed during the lab sessions and signed off by either module instructor or teaching assistant. If you are not able to complete the tasks during the lab sessions, then you should do them at home and have them ready to be marked off by the deadline. **All lab tasks also must be uploaded to Blackboard before the deadline stated on each assignment sheet.**

If there is any report or dissertation, it must be written and submitted in **PDF** format. Source codes must be organised and formatted neatly, sufficient and clear comments are very welcome and necessary for markers to assess your work. Submissions and feedback will be done via Blackboard-Tunitin system. Plagiarism will not be tolerated. Zip all your files with the following naming convention for submission:

- [Student Number]-[Last Name][First Initial]-[Assignment][Number].zip
- For example: *123456-DengJ-Assignment2.zip*

## CSCM35/CSLM35 Assignment 2 Complete by 12th/Mar/2020

This assignment is about getting familiar with an important supervised classification algorithm, Naïve Bayes classifier that was taught in the module.

### Iris Flower Dataset

The Iris flower data set or Fisher's Iris dataset is a multivariate dataset introduced by the British statistician and biologist Ronald Fisher in his 1936 paper. The dataset contains a set of 150 records under five attributes - petal length, petal width, sepal length, sepal width and species. The data set consists of 50 samples from each of three species of Iris (Iris setosa, Iris virginica and Iris versicolor), and four features were measured from each sample: the length and the width of the sepals and petals, in centimeters. A scatter plot of Fisher's Iris dataset is illustrated in Fig 1. The dataset is provided with *Python Scikit-Learn* package. More details on the data structure of the dataset and how to load using *Python* can be found from [https://scikit-learn.org/stable/auto\\_examples/datasets/plot\\_iris\\_dataset.html](https://scikit-learn.org/stable/auto_examples/datasets/plot_iris_dataset.html). In order to perform accuracy evaluation, you should split the dataset into two subsets for training (40 samples per category) and testing (10 samples per category). These two subsets should not have any overlap.

### □ Task – Classification with Naïve Bayes (10 marks)

Naïve Bayes classifiers are a family of simple “probabilistic classifier” based on applying Bayes' theorem with strong (naïve) independence assumptions between the features. In our case, we assume that the length and the width of the sepals and petals are independence random variables. The details of Naïve Bayes algorithm can be found at Section 8.3 of our recommended textbook (*Data Mining Concepts and Techniques, 3rd Edition*).

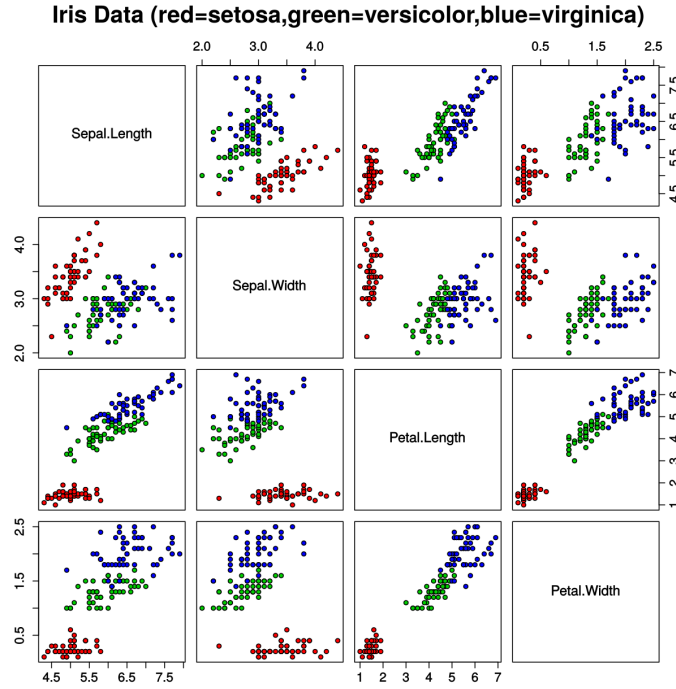


Figure 1: Scatterplot of Fisher's Iris dataset.

- **Cross Validation:** Split the dataset into two subsets, one for training (40 samples per category) and one for testing (10 samples per category). **(1/10 mark)**
- **Likelihood  $P(X_i|C_j)$ :** For individual feature of each category, construct a probabilistic model using Gaussian distribution, where you need to evaluate the mean and standard deviation of each feature across the samples which belong to the category.  $X_i$  represents a feature and  $C_j$  represents a category. **(3/10 marks)**
- **Priori  $P(C_j)$ :** Priori probabilities of individual categories can be evaluated using the sample population. **(1/10 marks)**
- **Posterior  $P(C_j|X)$ :** Given a testing sample  $X$ , you can calculate the posterior probability  $P(C_j|X)$  of a certain category  $C_j$  using Bayes' theorem and independence assumption of random variables (features). The category with highest posterior probability is voted as the prediction for given testing sample  $X$ . **(4/10 marks)**
- **Evaluation Confusion Matrix:** Calculate the prediction accuracy. **(1/10 mark)**
- Please note that the prediction accuracy is **NOT** important and it will **NOT** be considered for evaluate your work.
- Please submit your *Python Jupyter* code containing above calculations and describe your steps concisely using either markdown or comment whenever necessary.