# PRE-PROCESSING OF THE DATASET

## LAB 1

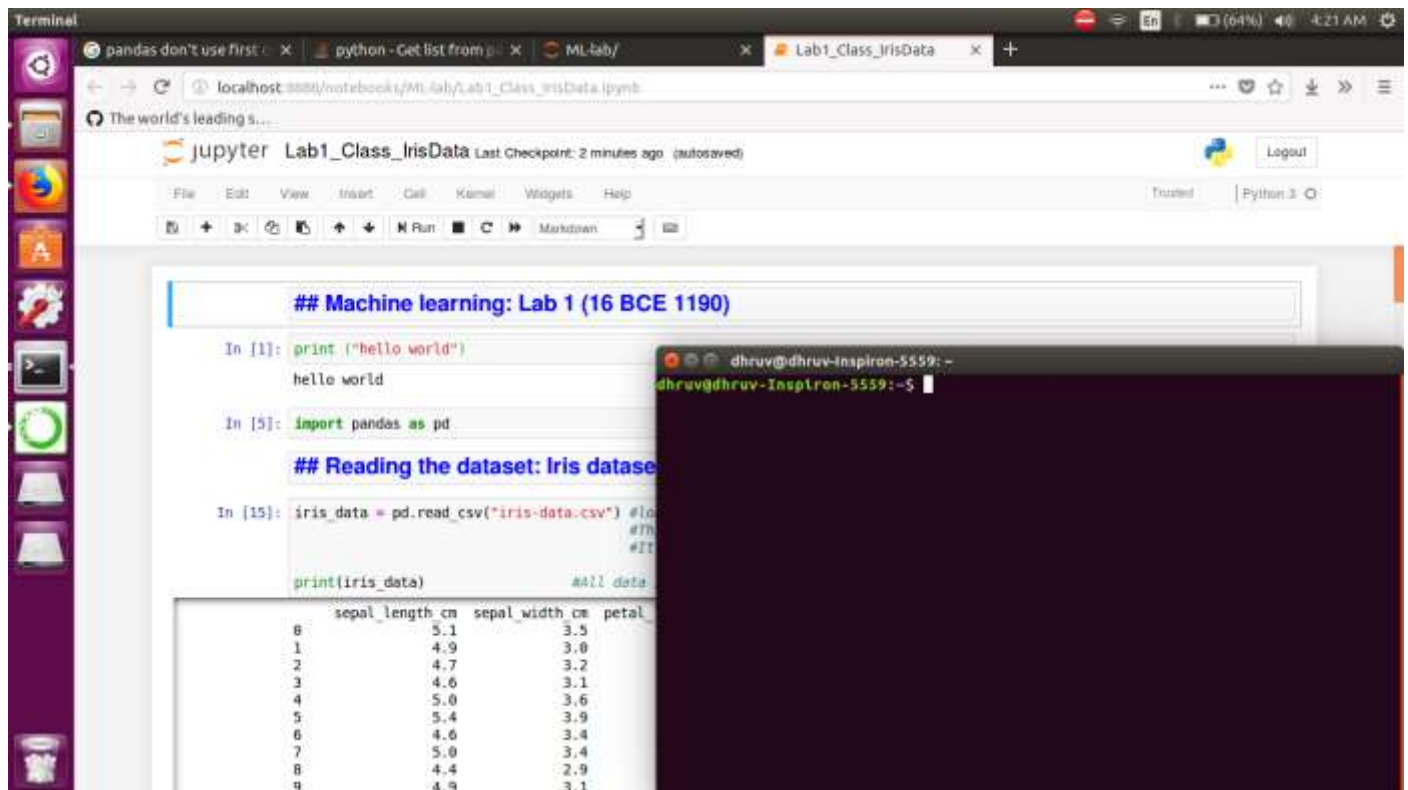**DATASET CHOSEN:**  Absenteeism at work (https://archive.ics.uci.edu/ml/datasets/Absenteeism+at+work)

The first step to implementing machine learning on any dataset is to **understand and clean the dataset** – in other words, perform pre-processing. In the process of understanding the dataset, we might find some faults in the data (**inconsistent measures, missing values**) – which will affect the algorithm implementation.

My first step was to import the pandas library and **read the ".csv" file using "read_csv"**. On reading, all the attribute values for a row went into the same column. This was found using the "shape" attribute of the dataframe. I observed that the delimiter used in the csv file was semi-colon(;) and not comma(,).  Thus, all values were incorrectly going into a single column. To correct this import of csv data, I added a " sep = ';' " argument in the read_csv command. Now when the first 5 rows of "absent_data" (dataframe) were printed using absent_data.head(), the columns and rows were correctly printed. I used the "describe", "shape", "index" and "columns" attributes to understand the dataset.
Next, was an important step – **to identify and correct null valued rows in the dataset**. This was done using "isnull().sum()" on the absent_data dataframe. Sequentially, all the columns with number of null valued rows were printed. It was noted that there were no null values in the dataset.

We know that machine learning algorithms are basically mathematical operations on the dataset. Hence, **math cannot be done on categorical (string/character) valued attributes.** We must convert them into numerical values. In this dataset, all values were already converted into numerical values by the dataset donator. But *for the sake of learning, 'Day of the week' attribute was changed from numerical to categorical and back to numerical value.* I tried using LabelEncoder() method of sklearn, but was unsuccessful. The encoder assigned numerical values to the days uniquely but randomly. This would change the meaning of the digits representing days in the original dataset.
Next, the **data** in the dataset **was scaled using the sklearn's StandardScaler.** Scaling of variables needs to be done so that no attribute gets higher/lower weight just because of the values it stores. Also, using the scaled values, computation of the algorithm is faster and so is its convergence.
Lastly, using the **seaborn** and **matplotlib** libraries, a **correlation plot was plotted between the various attributes of the dataset.** This helps us to know if there is a strong positive/negative correlation between two variables. Although the plot did not reveal any surprising insight about the dataset, it does lead us to a starting point for future analysis on the dataset.

**SCREENSHOTS**

**Iris-dataset**



**Chosen dataset – Absenteeism at work**