



**COMSATS University Islamabad**  
**Department of Electrical Engineering (Wah Campus)**  
**Artificial Intelligence (EEE-462) Lab Manual**

## **LAB # 04 Introduction to IRIS data set and scatter plots.**

### **Plotting for Exploratory data analysis (EDA)**

#### **Basic Terminology**

- What is EDA?
- Data-point/vector/Observation
- Data-set.
- Feature/Variable/Input-variable/Independent-varibale
- Label/dependent-variable/Output-variable/Class/Class-label/Response label
- Vector: 2-D, 3-D, 4-D,.... n-D

Q. What is a 1-D vector: Scalar

### **Iris Flower dataset**

**Toy Dataset:** Iris Dataset: [[https://en.wikipedia.org/wiki/Iris\\_flower\\_data\\_set](https://en.wikipedia.org/wiki/Iris_flower_data_set)]

- A simple dataset to learn the basics.
- 3 flowers of Iris species. [see images on wikipedia link above]
- 1936 by Ronald Fisher.
- Petal and Sepal: [http://terpconnect.umd.edu/~petersd/666/html/iris\\_with\\_labels.jpg](http://terpconnect.umd.edu/~petersd/666/html/iris_with_labels.jpg)
- **Objective:** Classify a new flower as belonging to one of the 3 classes given the 4 features.
- Importance of domain knowledge.
- Why use petal and sepal dimensions as features?
- Why do we not use 'color' as a feature?

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np
```

Download iris.csv from <https://raw.githubusercontent.com/uiuc-cse/data-fa14/gh-pages/data/iris.csv>

*Load Iris.csv into a pandas data Frame.*

```
iris = pd.read_csv("iris.csv")
```

how many data-points and features?

```
print (iris.shape)
(150, 5)
```



**COMSATS University Islamabad**  
**Department of Electrical Engineering (Wah Campus)**  
**Artificial Intelligence (EEE-462) Lab Manual**

What are the column names in our dataset?

```
print (iris.columns)
```

```
Index(['sepal_length', 'sepal_width', 'petal_length', 'petal_width',  
      'species'],  
      dtype='object')
```

How many data points for each class are present?

How many flowers for each species are present?

```
iris["species"].value_counts()  
# balanced-dataset vs imbalanced datasets  
#Iris is a balanced dataset as the number of data points for every class is 50.
```

```
virginica      50  
versicolor    50  
setosa         50  
Name: species, dtype: int64
```

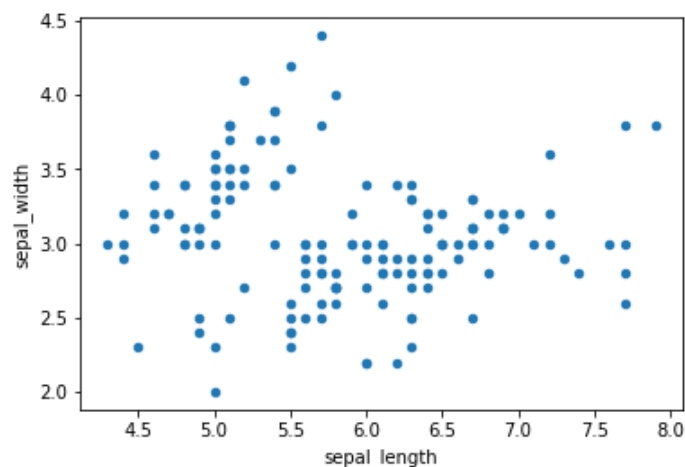
## 2-D Scatter Plot

ALWAYS understand the axis: labels and scale.

```
iris.plot(kind='scatter', x='sepal_length', y='sepal_width') ;  
plt.show()
```

cannot make much sense out of it.

What if we colour the points by their class-label/flower-type.



2-D Scatter plot with color-coding for each flower type/class.

Here 'sns' corresponds to seaborn.

```
sns.set_style("whitegrid");  
sns.FacetGrid(iris, hue="species", size=4) \  
    .map(plt.scatter, "sepal_length", "sepal_width") \  
    .add_legend();  
plt.show();
```

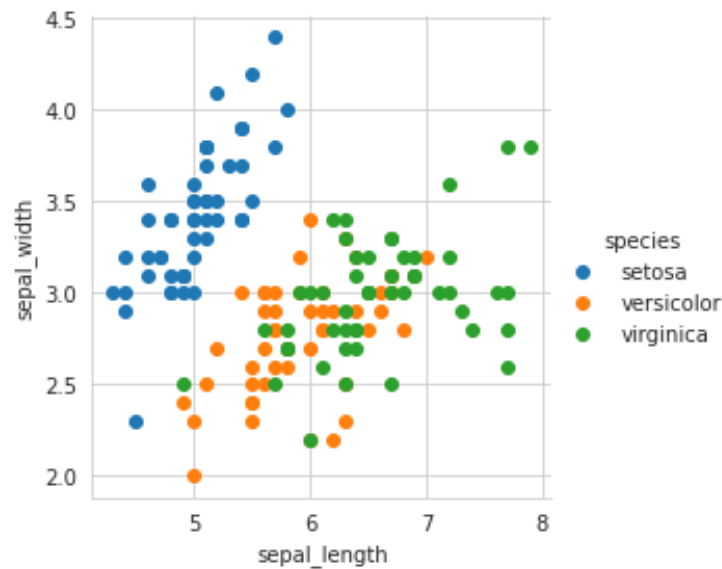


**COMSATS University Islamabad**  
**Department of Electrical Engineering (Wah Campus)**  
**Artificial Intelligence (EEE-462) Lab Manual**

Notice that the blue points can be easily separated from red and green by drawing a line. But red and green data points cannot be easily separated.

Can we draw multiple 2-D scatter plots for each combination of features?

How many combinations exist?  $4C2 = 6$ .



**Observation(s):**

Using sepal\_length and sepal\_width features, we can distinguish Setosa flowers from others.

Separating Versicolor from Virginica is much harder as they have considerable overlap.

**3D Scatter plot**

<https://plot.ly/pandas/3d-scatter-plots/>

Needs a lot of mouse interaction to interpret data.

What about 4-D, 5-D or n-D scatter plot?

**Pair-plot**

Pairwise scatter plot

Dis-advantages:

Can be used when number of features are high.

Cannot visualize higher dimensional patterns in 3-D and 4-D.

Only possible to view 2D patterns.

```
plt.close();  
sns.set_style("whitegrid");  
sns.pairplot(iris, hue="species", size=3);  
plt.show()
```



**COMSATS University Islamabad**  
**Department of Electrical Engineering (Wah Campus)**  
**Artificial Intelligence (EEE-462) Lab Manual**



### Observations

petal\_length and petal\_width are the most useful features to identify various flower types.

While Setosa can be easily identified (linearly separable), Versicolor and Virginica have some overlap (almost linearly separable).

We can find "lines" and "if-else" conditions to build a simple model to classify the flower types.

### Histogram, PDF, CDF

What about 1-D scatter plot using just one feature?

1-D scatter plot of petal-length



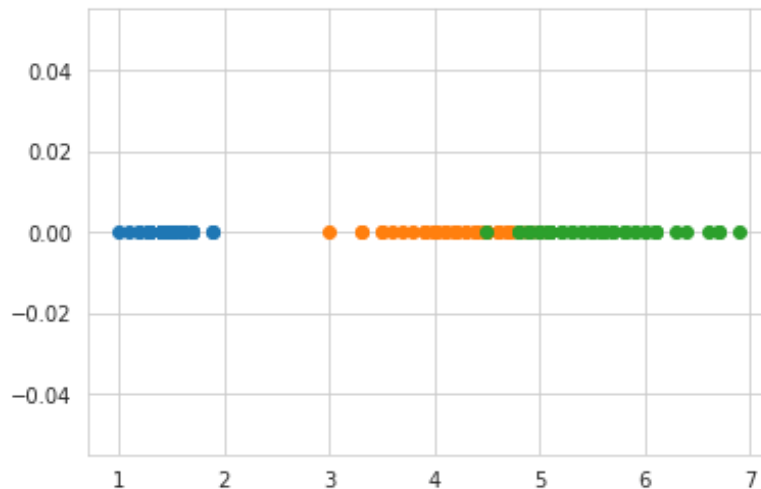
**COMSATS University Islamabad**  
**Department of Electrical Engineering (Wah Campus)**  
**Artificial Intelligence (EEE-462) Lab Manual**

```
import numpy as np
iris_setosa = iris.loc[iris["species"] == "setosa"];
iris_virginica = iris.loc[iris["species"] == "virginica"];
iris_versicolor = iris.loc[iris["species"] == "versicolor"];
#print(iris_setosa["petal_length"])
plt.plot(iris_setosa["petal_length"], np.zeros_like(iris_setosa['petal_length']), 'o')
plt.plot(iris_versicolor["petal_length"],
np.zeros_like(iris_versicolor['petal_length']), 'o')
plt.plot(iris_virginica["petal_length"],
np.zeros_like(iris_virginica['petal_length']), 'o')

plt.show()
```

Disadvantages of 1-D scatter plot: Very hard to make sense as points are overlapping a lot.

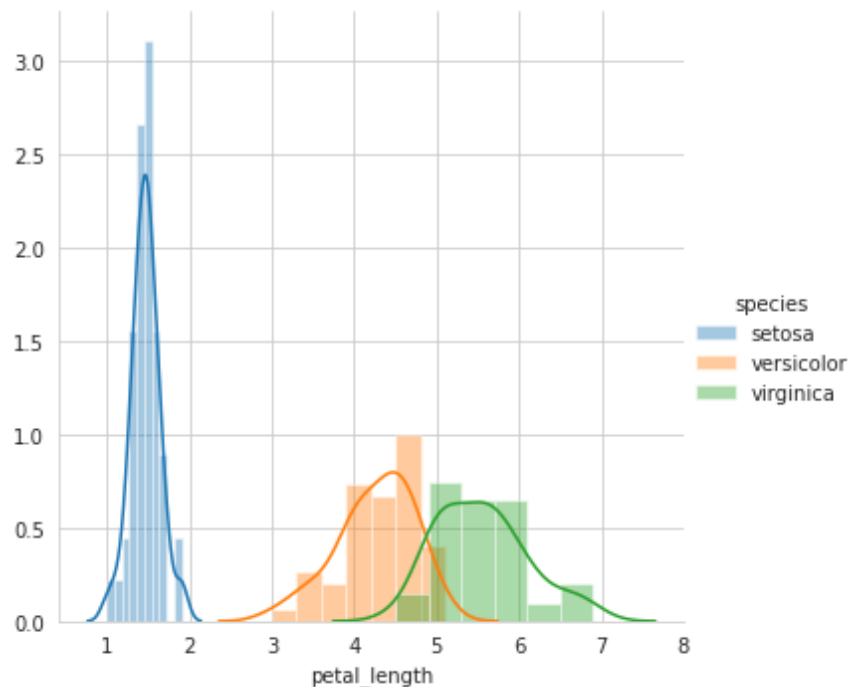
Are there better ways of visualizing 1-D scatter plots?



```
sns.FacetGrid(iris, hue="species", size=5) \
    .map(sns.distplot, "petal_length") \
    .add_legend();
plt.show();
```



**COMSATS University Islamabad**  
**Department of Electrical Engineering (Wah Campus)**  
**Artificial Intelligence (EEE-462) Lab Manual**



```
sns.FacetGrid(iris, hue="species", size=5) \
    .map(sns.distplot, "petal_width") \
    .add_legend();
plt.show();
```

```
sns.FacetGrid(iris, hue="species", size=5) \
    .map(sns.distplot, "sepal_length") \
    .add_legend();
plt.show();
```

```
sns.FacetGrid(iris, hue="species", size=5) \
    .map(sns.distplot, "sepal_width") \
    .add_legend();
plt.show();
```

### Histograms and Probability Density Functions (PDF)

How to compute PDFs using counts/frequencies of data points in each window.

How window width effects the PDF plot.

Interpreting a PDF:

why is it called a density plot?

Why is it called a probability plot?

for each value of petal\_length, what does the value on y-axis mean?

Notice that we can write a simple if..else condition as if(petal\_length) < 2.5 then flower type is setosa.

Using just one feature, we can build a simple "model" using if..else... statements.



**COMSATS University Islamabad**  
**Department of Electrical Engineering (Wah Campus)**  
**Artificial Intelligence (EEE-462) Lab Manual**

Disadv of PDF: Can we say what percentage of versicolor points have a petal\_length of less than 5?

Do some of these plots look like a bell-curve you studied in under-grad?

Gaussian/Normal distribution.

What is "normal" about normal distribution?

e.g: Heights of male students in a class.

One of the most frequent distributions in nature.

### Need for Cumulative Distribution Function (CDF)

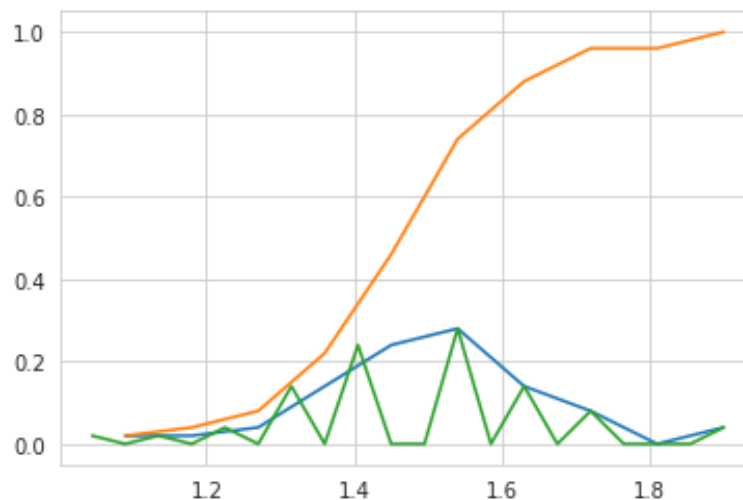
We can visually see what percentage of versicolor flowers have a petal\_length of less than 5?

How to construct a CDF?

How to read a CDF?

Plot CDF of petal\_length

```
counts, bin_edges = np.histogram(iris_setosa['petal_length'], bins=10,  
                                density = True)  
  
pdf = counts/(sum(counts))  
print(pdf);  
print(bin_edges);  
cdf = np.cumsum(pdf)  
plt.plot(bin_edges[1:],pdf);  
plt.plot(bin_edges[1:], cdf)  
  
counts, bin_edges = np.histogram(iris_setosa['petal_length'], bins=20,  
                                density = True)  
  
pdf = counts/(sum(counts))  
plt.plot(bin_edges[1:],pdf);  
  
plt.show();
```





**COMSATS University Islamabad**  
**Department of Electrical Engineering (Wah Campus)**  
**Artificial Intelligence (EEE-462) Lab Manual**

Need for Cumulative Distribution Function (CDF)

We can visually see what percentage of versicolor flowers have a petal\_length of less than 1.6?

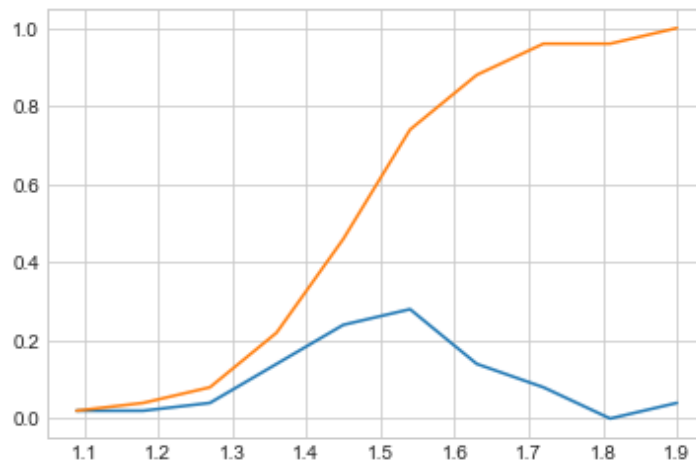
How to construct a CDF?

How to read a CDF?

Plot CDF of petal\_length

```
counts, bin_edges = np.histogram(iris_setosa['petal_length'], bins=10,
                                  density = True)

pdf = counts/(sum(counts))
print(pdf);
print(bin_edges)
#compute CDF
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:], cdf)
plt.show();
```



Plots of CDF of petal\_length for various types of flowers.

Misclassification error if you use petal\_length only.

```
counts, bin_edges = np.histogram(iris_setosa['petal_length'], bins=10,density = True)
pdf = counts/(sum(counts))
print(pdf);
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:], cdf)

# virginica
counts, bin_edges = np.histogram(iris_virginica['petal_length'], bins=10, density =
True)
pdf = counts/(sum(counts))
print(pdf);
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:], cdf)

#versicolor
```





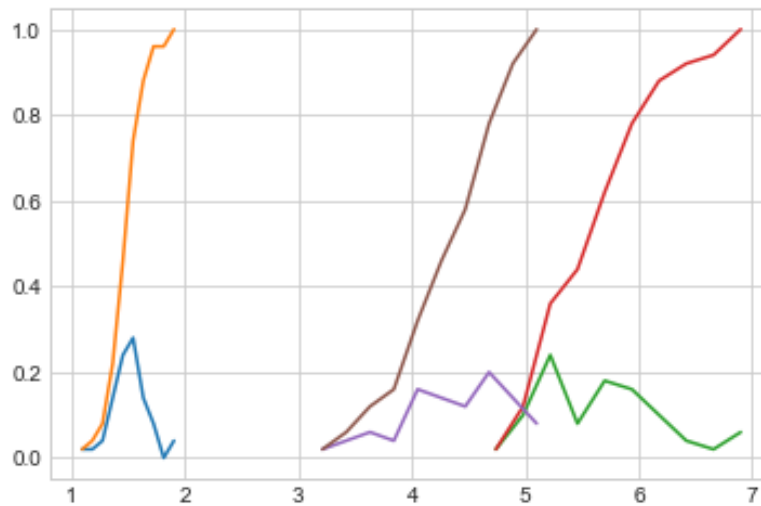
# COMSATS University Islamabad

## Department of Electrical Engineering (Wah Campus)

### Artificial Intelligence (EEE-462) Lab Manual

```
counts, bin_edges = np.histogram(iris_versicolor['petal_length'], bins=10, density =
True)
pdf = counts/(sum(counts))
print(pdf);
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:], cdf)

plt.show();
```



## Mean, Variance and Std-dev

Mean, Variance, Std-deviation,

```
print("Means:")
print(np.mean(iris_setosa["petal_length"]))
```

Mean with an outlier.

```
print(np.mean(np.append(iris_setosa["petal_length"],50)));
print(np.mean(iris_virginica["petal_length"]))
print(np.mean(iris_versicolor["petal_length"]))

print("\nStd-dev:");
print(np.std(iris_setosa["petal_length"]))
print(np.std(iris_virginica["petal_length"]))
print(np.std(iris_versicolor["petal_length"]))
```

## Median, Percentile, Quantile, IQR, MAD

```
#Median, Quantiles, Percentiles, IQR.
print("\nMedians:")
print(np.median(iris_setosa["petal_length"]))

#Median with an outlier
print(np.median(np.append(iris_setosa["petal_length"],50)));
print(np.median(iris_virginica["petal_length"]))
print(np.median(iris_versicolor["petal_length"]))
```



# COMSATS University Islamabad

## Department of Electrical Engineering (Wah Campus)

### Artificial Intelligence (EEE-462) Lab Manual

```
print("\nQuantiles:")
print(np.percentile(iris_setosa["petal_length"], np.arange(0, 100, 25)))
print(np.percentile(iris_virginica["petal_length"], np.arange(0, 100, 25)))
print(np.percentile(iris_versicolor["petal_length"], np.arange(0, 100, 25)))

print("\n90th Percentiles:")
print(np.percentile(iris_setosa["petal_length"], 90))
print(np.percentile(iris_virginica["petal_length"], 90))
print(np.percentile(iris_versicolor["petal_length"], 90))

from statsmodels import robust

print("\nMedian Absolute Deviation")
print(robust.mad(iris_setosa["petal_length"]))
print(robust.mad(iris_virginica["petal_length"]))
print(robust.mad(iris_versicolor["petal_length"]))
```

### Box plot and Whiskers

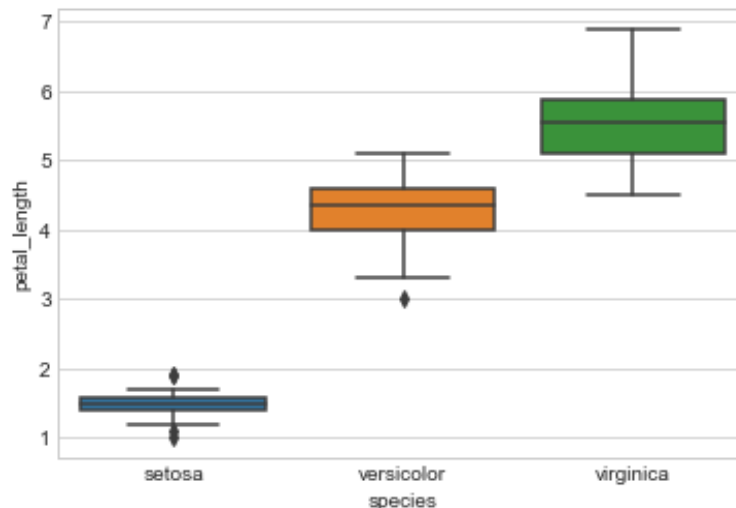
- Box-plot with whiskers: another method of visualizing the 1-D scatter plot more intuitively.
- The Concept of median, percentile, quantile.
- How to draw whiskers: [no standard way] Could use min and max or use other complex statistical techniques.
- IQR like idea.

NOTE: IN the plot below, a technique call inter-quartile range is used in plotting the whiskers.

Whiskers in the plot below donot correposnd to the min and max values.

Box-plot can be visualized as a PDF on the side-ways.

```
sns.boxplot(x='species', y='petal_length', data=iris)
plt.show()
```





# COMSATS University Islamabad

## Department of Electrical Engineering (Wah Campus)

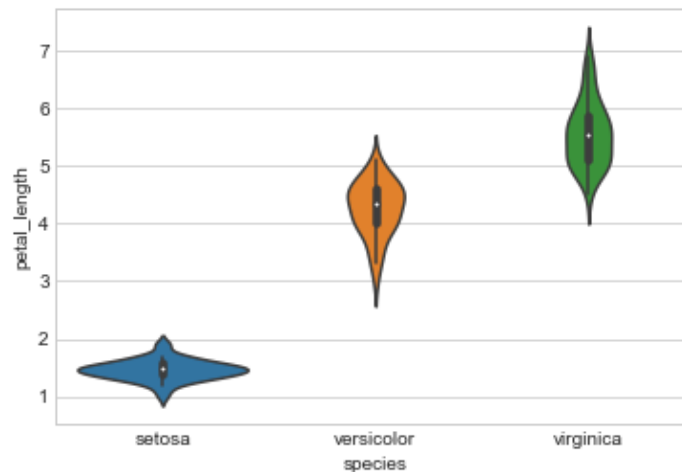
### Artificial Intelligence (EEE-462) Lab Manual

#### Violin plots

A violin plot combines the benefits of the previous two plots and simplifies them

Denser regions of the data are fatter, and sparser ones thinner in a violin plot

```
sns.violinplot(x="species", y="petal_length", data=iris, size=8)
plt.show()
```



#### Summarizing plots in english

Explain your findings/conclusions in plain english

Never forget your objective (the problem you are solving).

Perform all of your EDA aligned with your objectives.

#### Univariate, bivariate and multivariate analysis.

```
Def: Univariate, Bivariate and Multivariate analysis.
File "<ipython-input-20-f25211abae88>", line 3
    Def: Univariate, Bivariate and Multivariate analysis.
        ^
SyntaxError: invalid syntax
```

#### Multivariate probability density, contour plot.

```
#2D Density plot, contours-plot
sns.jointplot(x="petal_length", y="petal_width", data=iris_setosa, kind="kde");
plt.show();
```

#### Lab Tasks:

Download Haberman Cancer Survival dataset from Kaggle. You may have to create a Kaggle account to download data. (<https://www.kaggle.com/gilsousa/habermans-survival-data-set>)

Perform a similar analysis as above on this dataset with the following sections:



# **COMSATS University Islamabad**

## **Department of Electrical Engineering (Wah Campus)**

### **Artificial Intelligence (EEE-462) Lab Manual**

High level statistics of the dataset: number of points, number of features, number of classes, data-points per class.

Explain our objective.

Perform Univariate analysis(PDF, CDF, Boxplot, Violin plots) to understand which features are useful towards classification.

Perform Bi-variate analysis (scatter plots, pair-plots) to see if combinations of features are useful in classification.

Write your observations in english as crisply and unambiguously as possible. Always quantify your results.

```
iris_virginica_SW = iris_virginica.iloc[:,1]
iris_versicolor_SW = iris_versicolor.iloc[:,1]

from scipy import stats
stats.ks_2samp(iris_virginica_SW, iris_versicolor_SW)

x = stats.norm.rvs(loc=0.2, size=10)
stats.kstest(x, 'norm')

x = stats.norm.rvs(loc=0.2, size=100)
stats.kstest(x, 'norm')

x = stats.norm.rvs(loc=0.2, size=1000)
stats.kstest(x, 'norm')
```