

3. Plotting for Exploratory data analysis (EDA)

(3.1) Basic Terminology

- What is EDA?
 - Data-point/vector/Observation
 - Data-set
 - Feature/Attribute/input-variable/Dependent-variable
 - Label/Independent-variable/Output-varible/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

- 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://www.biorxiv.org/content/10.1101/066746v1.full.pdf>
- 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?

```
In [3]: 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/ucuc-cse/data-fair4/gh-pages/data/iris.csv'''
iris = pd.read_csv("E:/applied/iris.csv")

In [4]: # (Q) How many data-points and features?
print(iris.shape)
(150, 5)

In [5]: # (Q) What are the column names in our dataset?
print(iris.columns)
Index(['sepal_length', 'sepal_width', 'petal_length', 'petal_width',
       'species'],
      dtype='object')

In [6]: # (Q) How many data points for each class are present?
# (Q) 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.

Out[6]: virginica      50
setosa      50
versicolour  50
Name: species, dtype: int64
```

(3.2) 2-D Scatter Plot

```
In [7]: # 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 color the points by their class-label/flower-type.

In [8]: # 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() \
    .show()

# 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):

1. Using sepal_length and sepal_width features, we can distinguish Setosa flowers from others.
2. Separating Versicolour from Virginica is much harder as they have considerable overlap.

3D Scatter plot

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

Needs a lot to mouse interaction to interpret data.

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

(3.3) Pair-plot

```
In [9]: # pairwise scatter plot: Pair-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()
# NOTE: the diagonal elements are PDFs for each feature. PDFs are explained below.
```

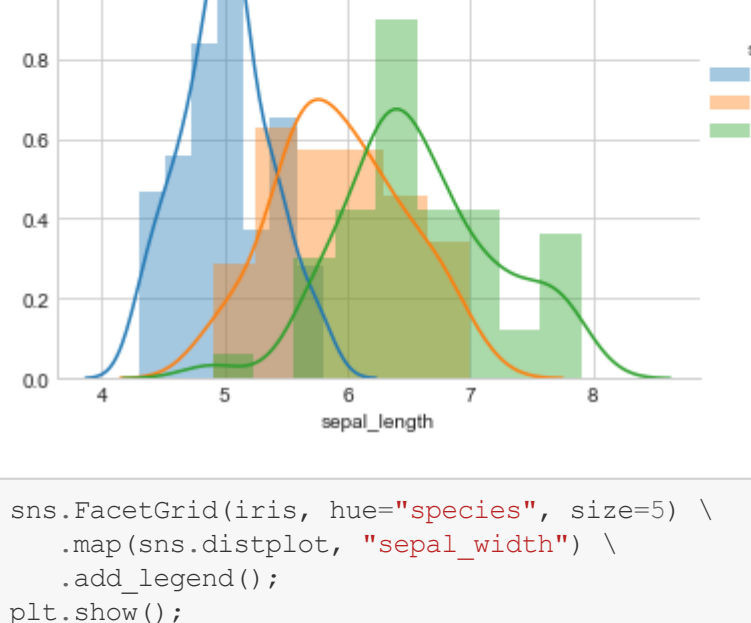


Observations

1. petal_length and petal_width are the most useful features to identify various flower types.
2. While Setosa can be easily identified (linearly separable), Versicolour and Versicolour have some overlap (almost linearly separable).
3. We can find "and" and "if-else" conditions to build a simple model to classify the flower types.

(3.4) Histogram, PDF, CDF

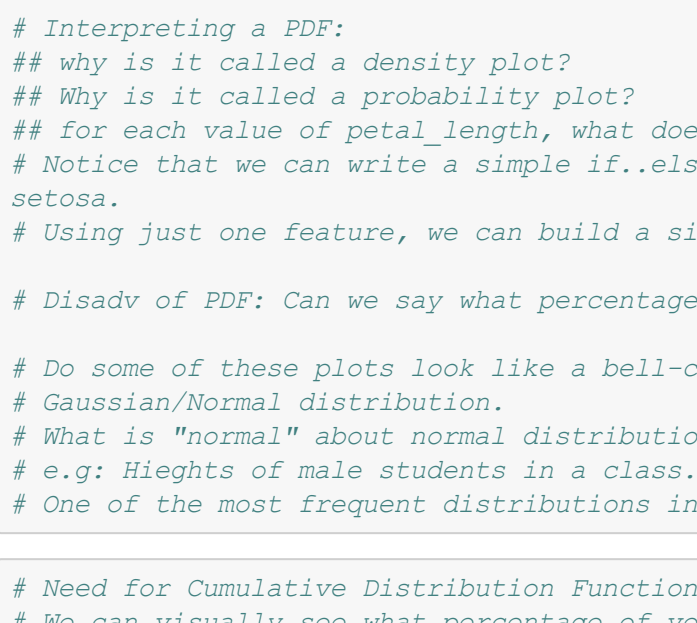
```
In [10]: # What about 1-D scatter plot using just one feature?
# 1-D scatter plot of petal-length
import numpy as np
iris_setosa = iris.loc[iris['species'] == "setosa"]
iris_virginica = iris.loc[iris['species'] == "virginica"]
iris_versicolour = iris.loc[iris['species'] == "versicolour"]
# Print(iris_setosa["petal_length"])
plt.plot(iris_setosa["petal_length"], np.zeros_like(iris_setosa["petal_length"]), 'o')
plt.plot(iris_versicolour["petal_length"], np.zeros_like(iris_versicolour["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?
```



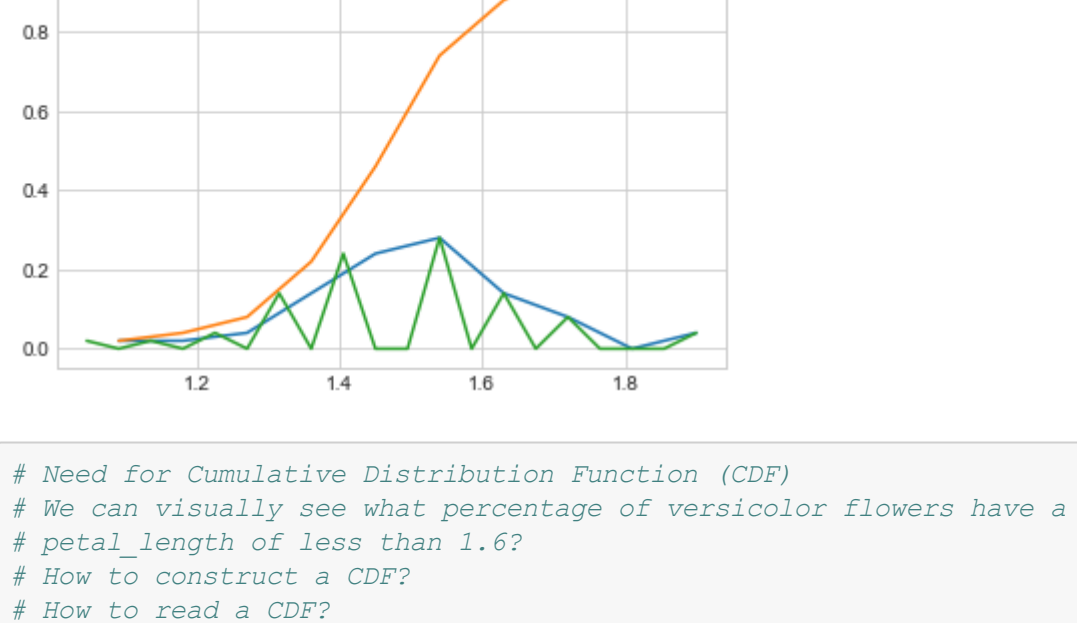
```
In [11]: sns.FacetGrid(iris, hue='species', size=3) \
    .map(sns.distplot, "petal_length") \
    .add_legend() \
    .show()

In [12]: sns.FacetGrid(iris, hue='species', size=3) \
    .map(sns.distplot, "sepal_width") \
    .add_legend() \
    .show()

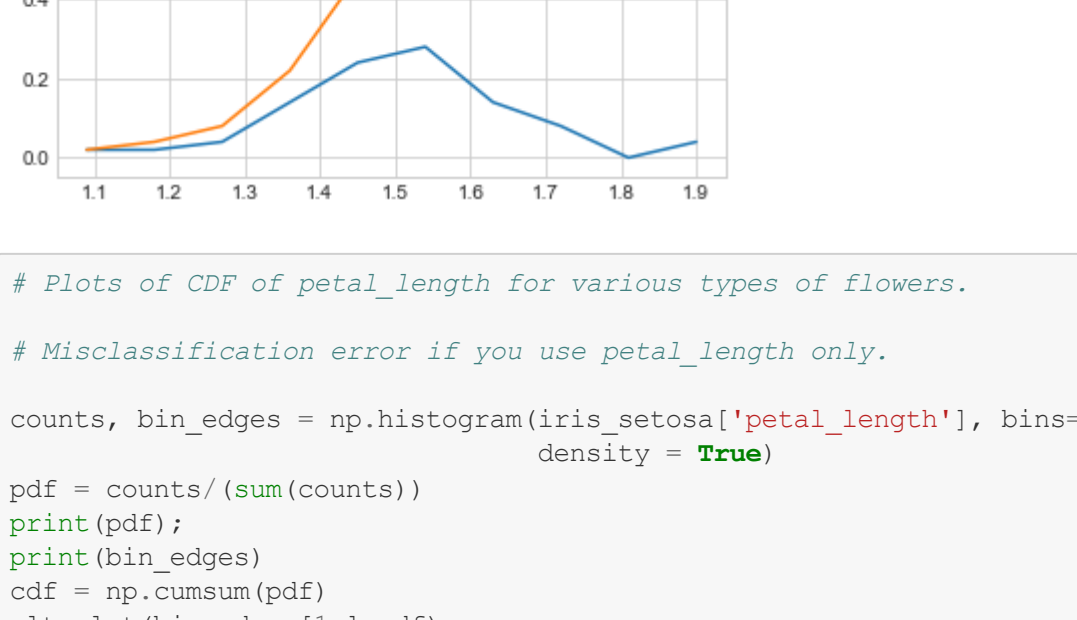
In [13]: sns.FacetGrid(iris, hue='species', size=3) \
    .map(sns.distplot, "sepal_length") \
    .add_legend() \
    .show()
```



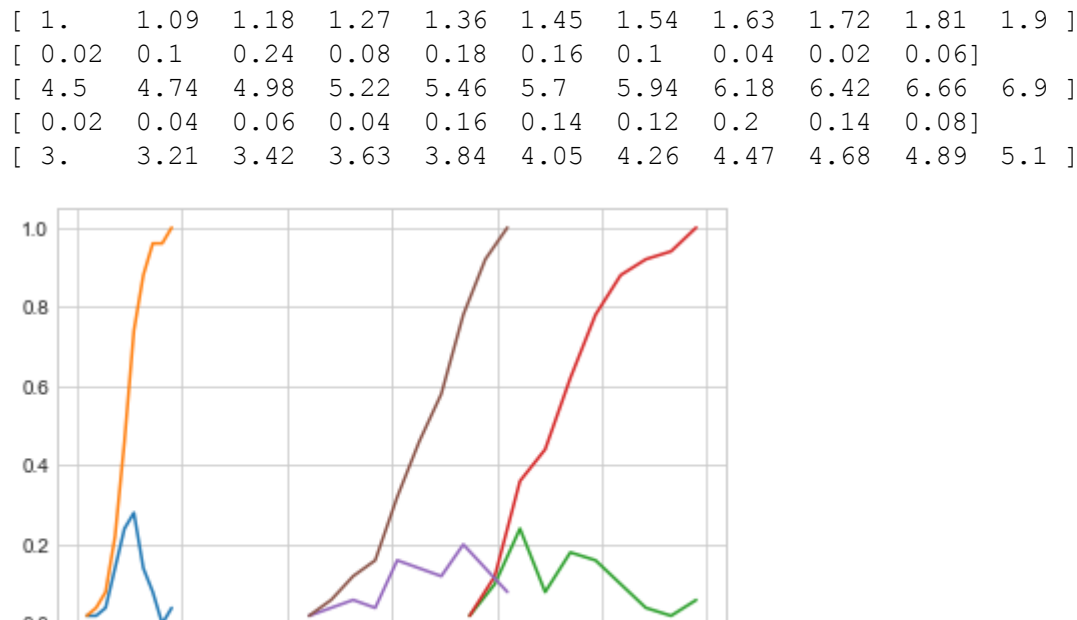
```
In [14]: # Need for Cumulative Distribution Function (CDF)
# We can visually see what percentage of versicolour flowers have a
# petal_length of less than 5?
# How to construct a CDF?
# How to read a CDF?
# Plot CDF of petal_length
counts, bin_edges = np.histogram(iris_setosa["petal_length"], bins=10,
                                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)
```



```
In [15]: # Need for Cumulative Distribution Function (CDF)
# We can visually see what percentage of versicolour 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()
```

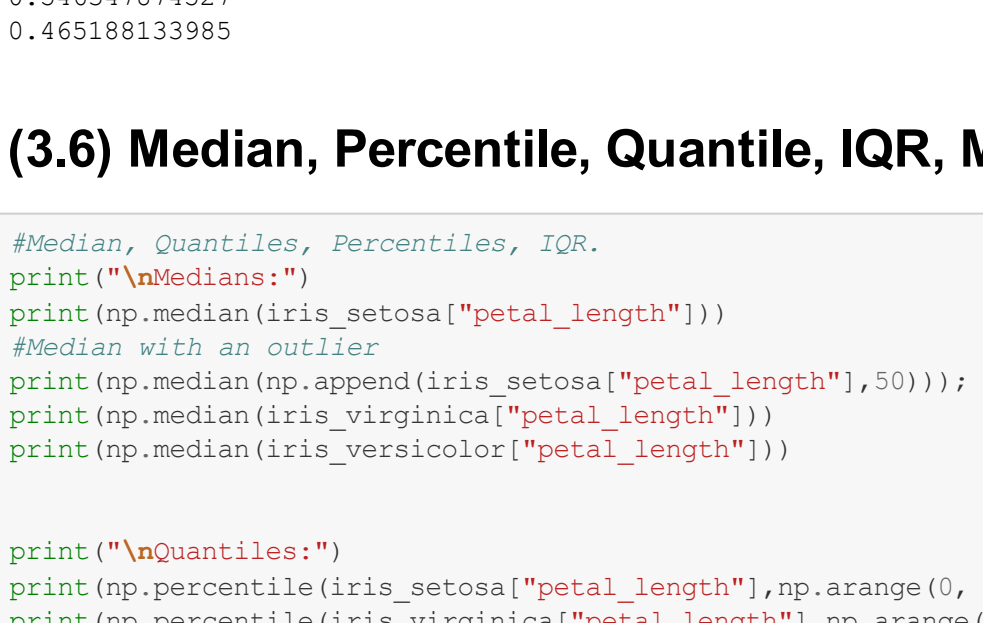


```
In [16]: # Plots of CDF of petal_length for various types of flowers.
# Misclassification error if you use petal_length only.
counts, bin_edges = np.histogram(iris_setosa["petal_length"], bins=10,
                                density = True)
pdf = counts/(sum(counts))
print(pdf)
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:], cdf)
```



```
In [17]: # 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)

# versicolour
counts, bin_edges = np.histogram(iris_versicolour["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()
```



(3.5) Mean, Variance and Std-dev

```
In [18]: # 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_versicolour["petal_length"]))

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

Median, Percentile, Quantile, IQR, MAD

```
In [19]: #Median, Quantiles, Percentile, IQR, MAD
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_versicolour["petal_length"]))

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_versicolour["petal_length"], np.arange(0, 100, 25)))

print("\n100th Percentiles:")
print(np.percentile(iris_setosa["petal_length"], 90))
print(np.percentile(iris_virginica["petal_length"], 90))
print(np.percentile(iris_versicolour["petal_length"], 90))

from statsmodels import robust
print(robust.mad(iris_setosa["petal_length"]))
print(robust.mad(iris_virginica["petal_length"]))
print(robust.mad(iris_versicolour["petal_length"]))

Medians:
1.5
1.5
4.35
4.35

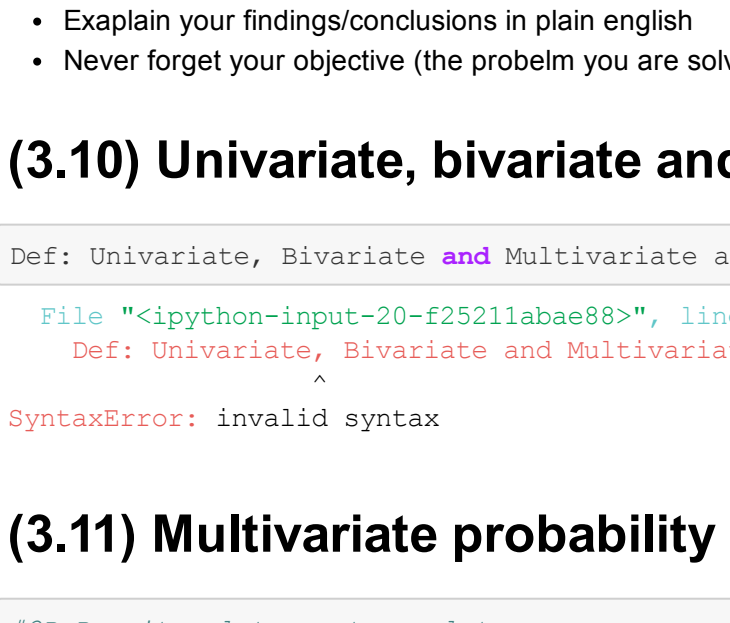
Quantiles:
[ 1.    1.4   1.5   1.575]
[ 1.65   5.1   5.55  5.875]
[ 3.    4.    4.35  4.6 ]

90th Percentiles:
1.7
6.31
4.8

Median Absolute Deviation
0.14246221831
0.667170998328
0.518910778477
```

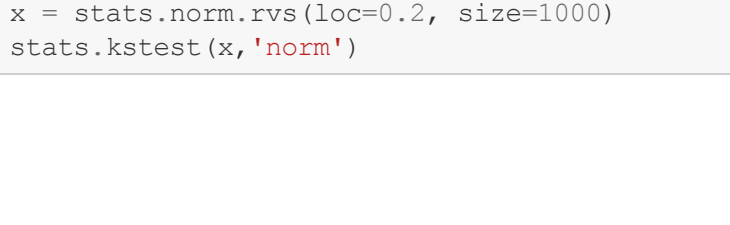
(3.7) Box plot and Whiskers

```
In [20]: #Box-plot with whiskers: another method of visualizing the 1-D scatter plot more intuitively.
# The Concept of median, percentile, quantile
# How to draw the box in the box-plot?
# How to draw whiskers: (no standard way) could use min and max or use other complex statistical test
# 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 do not correspond to the min and max values.
#Box-plot can be visualized as a PDF on the side-way.
sns.boxplot(x='species', y='petal_length', data=iris)
plt.show()
```



(3.8) Violin plots

```
In [21]: # 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()
```



(3.9) 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.

(3.10) Univariate, bivariate and multivariate analysis.

```
In [22]: Def: Univariate, Bivariate and Multivariate analysis.
File "C:\python-input-20-625211abae88>", line 3
Def: Univariate, Bivariate and Multivariate analysis.
>
SyntaxError: invalid syntax
```

(3.11) Multivariate probability density, contour plot.

```
In [23]: #2D Density plot, contour-plot
sns.jointplot(x='petal_length', y='petal_width', data=iris_setosa, kind='kde');
plt.show()
```

(3.12) Exercise:

1. Download Haberman Cancer Survival dataset from Kaggle. You may have to create a Kaggle account to download data (<https://www.kaggle.com/habermanhabermansurvival/data>)
2. Perform a similar analysis as above on this dataset with the following sections:
3. High level statistics of the dataset: number of points, number of features, number of classes, data-points per class.
4. Explain our objective.
5. Perform Univariate analysis(PDF, CDF, Boxplot, Violin plots) to understand which features are useful towards classification
6. Perform Bi-variate analysis (scatter plots, pair-plots) to see if combinations of features are useful in classification.
7. Write your observations in english as crisply and unambiguously as possible. Always quantify your results.

```
In [24]: iris_virginica_SW = iris_virginica.loc[:,1:]
iris_versicolour_SW = iris_versicolour.loc[:,1:]

In [25]: from scipy import stats
stats.kstest(iris_virginica_SW, iris_versicolour_SW)

In [26]: x = stats.norm.rvs(loc=0.2, size=100)
stats.kstest(x, 'norm')

In [27]: x = stats.norm.rvs(loc=0.2, size=1000)
stats.kstest(x, 'norm')
```